

Development and Use of Predictive Toxicology Applications

OpenTox Workshop

19 September 2010

Barry Hardy (Douglas Connect)

Rhodes, Greece

Collaborating Partners

In Silico Toxicology,
Switzerland

Douglas Connect,
Switzerland

Albert Ludwigs University
Freiburg, Germany

Ideaconsult,
Bulgaria

Istituto Superiore
di Sanità, Italy

Technical University
of Munich, Germany



National Technical
University of Athens,
Greece

Fraunhofer Institute
for Toxicology &
Experimental Medicine,
Germany

David Gallagher, UK

Institute of Biomedical
Chemistry of the Russian
Academy of Medical
Sciences, Russia

Seascope Learning &
JNU, India

OpenTox Advisory Board

- European Centre for the Validation of Alternative Methods
- Pharmatropé
- Bioclipse
- U.S. Environmental Protection Agency
- U.S. Food & Drug Administration
- Nestlé
- Roche
- AstraZeneca
- LHASA
- Leadscope
- University of North Carolina
- EC Environment Directorate General
- Organisation for Economic Cooperation & Development
- CADASTER
- Bayer Healthcare

Journal of Cheminformatics Publication

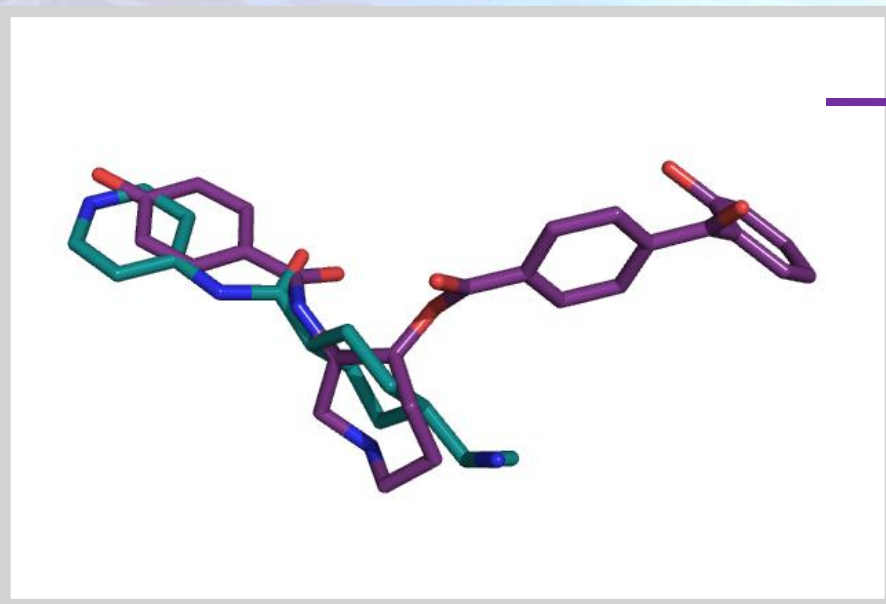
Collaborative development of predictive toxicology applications
Journal of Cheminformatics 2010, 2:7 doi:10.1186/1758-2946-2-7

Barry Hardy, Nicki Douglas, Christoph Helma, Micha Rautenberg, Nina Jeliaskova, Vedrin Jeliaskov, Ivelina Nikolova, Romualdo Benigni, OlgaTcheremenskaia, Stefan Kramer, Tobias Girschick, Fabian Buchwald, JoergWicker, Andreas Karwath, Martin Gutlein, Andreas Maunz, Haralambos Sarimveis, Georgia Melagraki, Antreas Afantitis, Pantelis Sopasakis, David Gallagher, Vladimir Poroikov, Dmitry Filimonov, Alexey Zakharov, Alexey Lagunin, Tatyana Gloriovova, Sergey Novikov, Natalia Skvortsova, Dmitry Druzhilovsky, Sunil Chawla, Indira Ghosh, Surajit Ray, Hitesh Patel and Sylvia Escher

Open Access publication available at
www.jcheminf.com/content/2/1/7

Collaborative Predictive Toxicology Challenge

Input Structure



VO

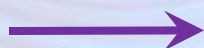


Out - Toxic or Not?

- ☐ LD50
- ☐ Liver Toxicity
- ☐ Secondary Metabolites
- ☐ Bioavailability
- ☐ Mutagenicity
- ☐ Carcogenicity
- ☐ Reproductive Toxicology
- ☐ Skin Irritation
- ☐ Aqua Toxicity
- ☐ Combined predictions for arrays of multiple end points



Driver



Increasing demands on industry to satisfy safety evaluation and risk assessment required by REACH legislation. (Over 140k cmpds registered).

Step 1: Search
 Select structure(s)

Step 2: Verify structure
 Verify structure

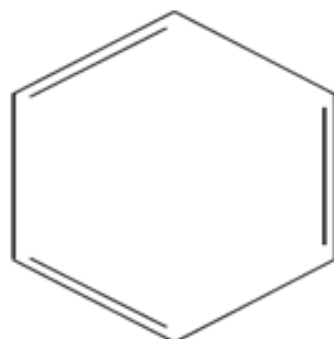
Step 3: Models
 Select prediction models

Step 4: Estimate
 Estimate

Step 5: Results
 Display results

This page lists your ToxPredict workflow results for the structure(s) you have selected and the model prediction(s) you have chosen to run. You could also retrieve the ToxPredict report in various other formats, e.g. [SDF](#), [CML](#), [SMI](#), [PDF](#), [CSV](#), [ARFF](#), [RDF/XML](#) or [RDF/N3](#).

Download as 



CAS RN
EINECS
IUPAC name
Synonym

71-43-2
 200-753-7
 benzene
 (6)annulene; benzine; Benzol; Benzolene;
 bicarburet of hydrogen; carbon oil; Coal naphtha;
 cyclohexatriene; mineral naphtha; motor benzol;
 nitration benzene; Phene; Phenyl hydride;
 pyrobenzol.

Synonym
Synonym
Synonym
Quality label

21742.0
 Benzene
 benzene
 OK

MolecularWeight  **MolecularWeight**

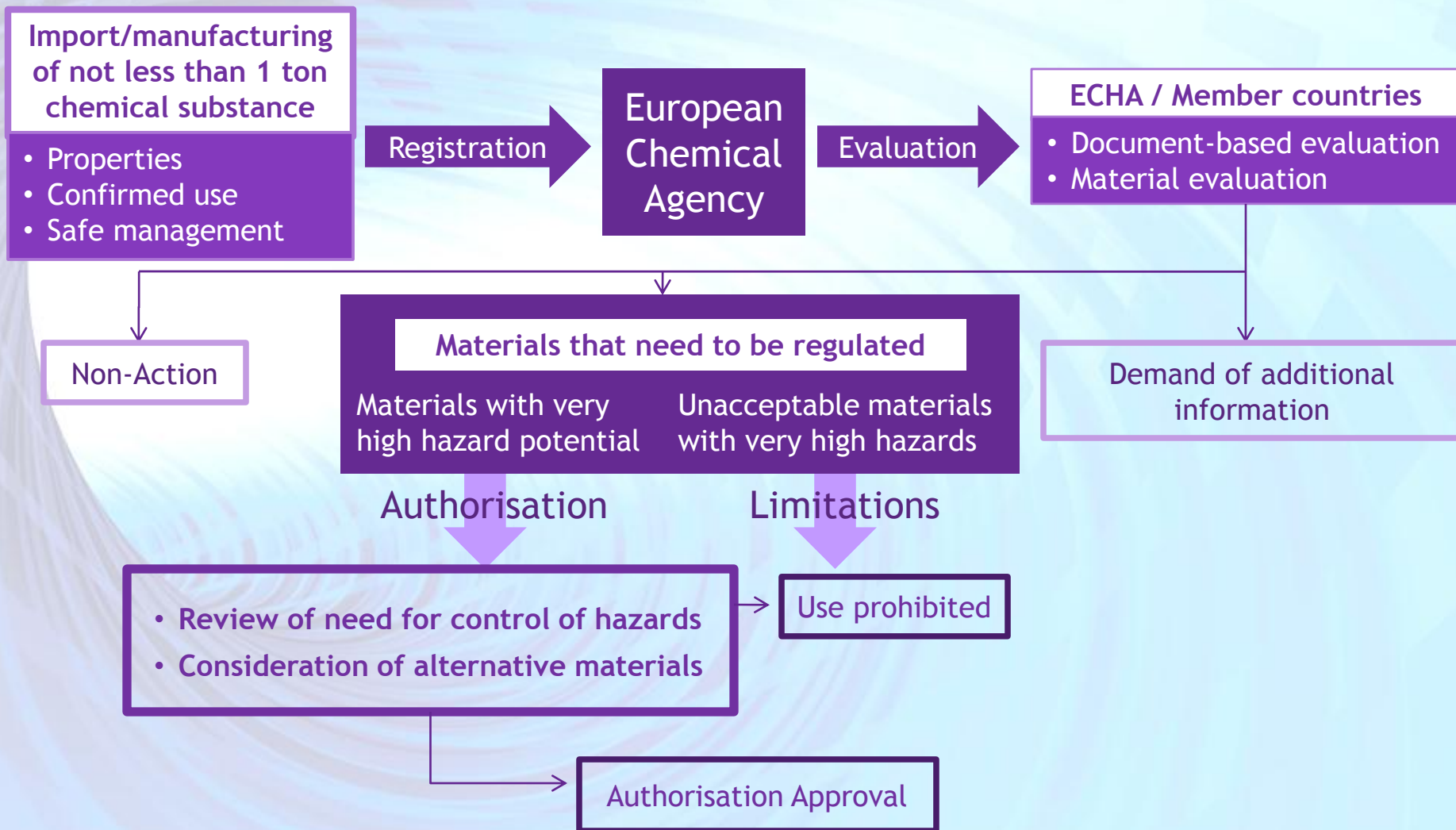
MW

78.1112

REACH



REACH Registration



Challenges to Integrated Resources & Applications

- Database silos
- Missing information
- Varying quality
- Hard to integrate data
- Hard to integrate models
- No common framework
- Lack of standards
- Lack of validation
- Complex subject
- Application difficult
- Lack of transparency
- Interdisciplinary collaboration

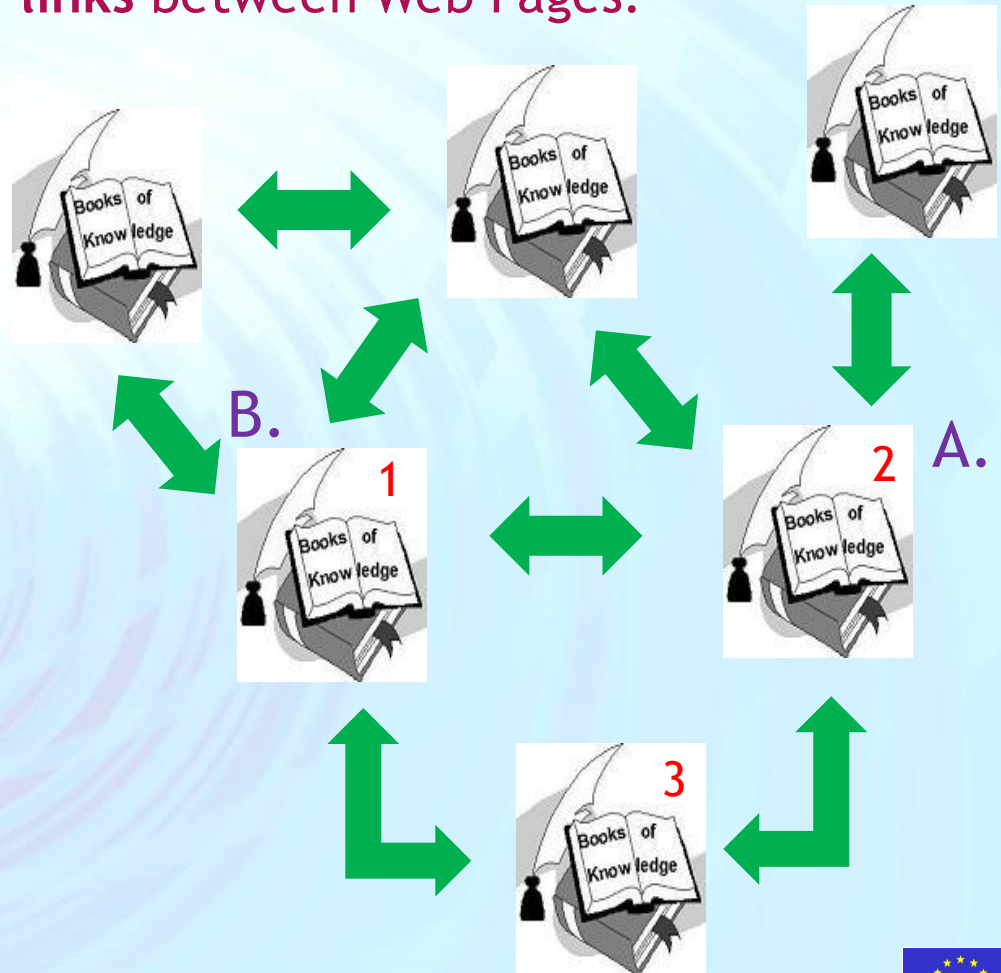
Value is in Linking

The key idea of Google's founders in creating their search engine:
There is useful knowledge in the **links** between Web Pages.

Page Ranking

A page is ranked higher in a search if:

- A. it has more connections to it than other pages
- B. the pages connecting to it have higher ranking themselves



Linked Data enables Knowledge Creation, Combination and Analysis

Linked Data is a term used to describe the exposing, sharing, and connecting of data on the Semantic Web using:

URIs a generic means to identify entities in the world

HTTP a simple yet universal mechanism for retrieving resources

RDF a generic graph-based data model with which to structure and link data

Linked Data needs:

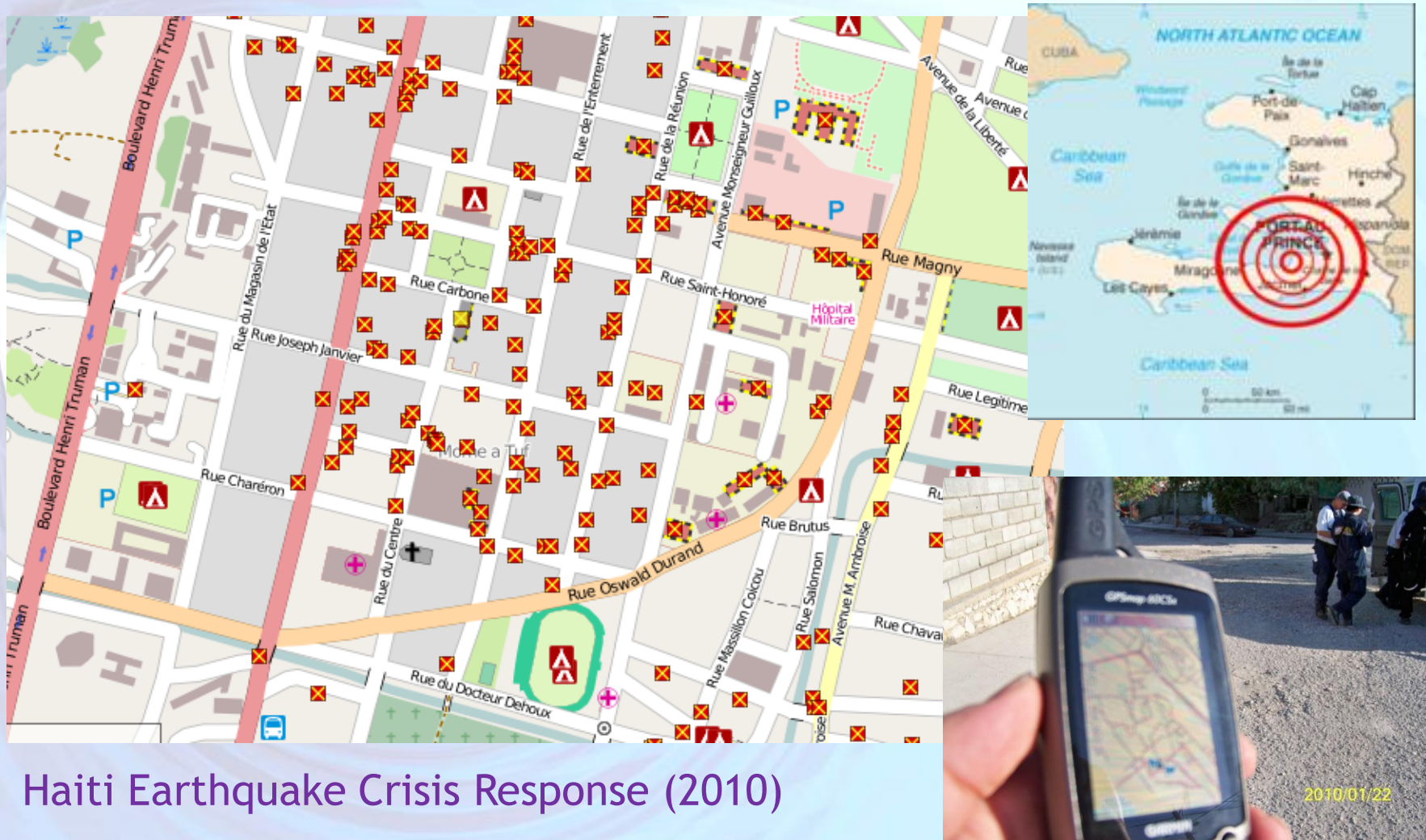
1. Provision of a **URI** that describes a Data Resource
2. Use of **HTTP** to retrieve useful data from the **URI**
3. A Data Format described with standardised semantics (so relationships are enabled) e.g. **RDF**
4. Data should provide links to other Data (through **URIs**)

Linked Data approach can also be applied to other resource types e.g., for algorithms or models as done in OpenTox...



DBpedia = Linked Data approach applied to Wikipedia

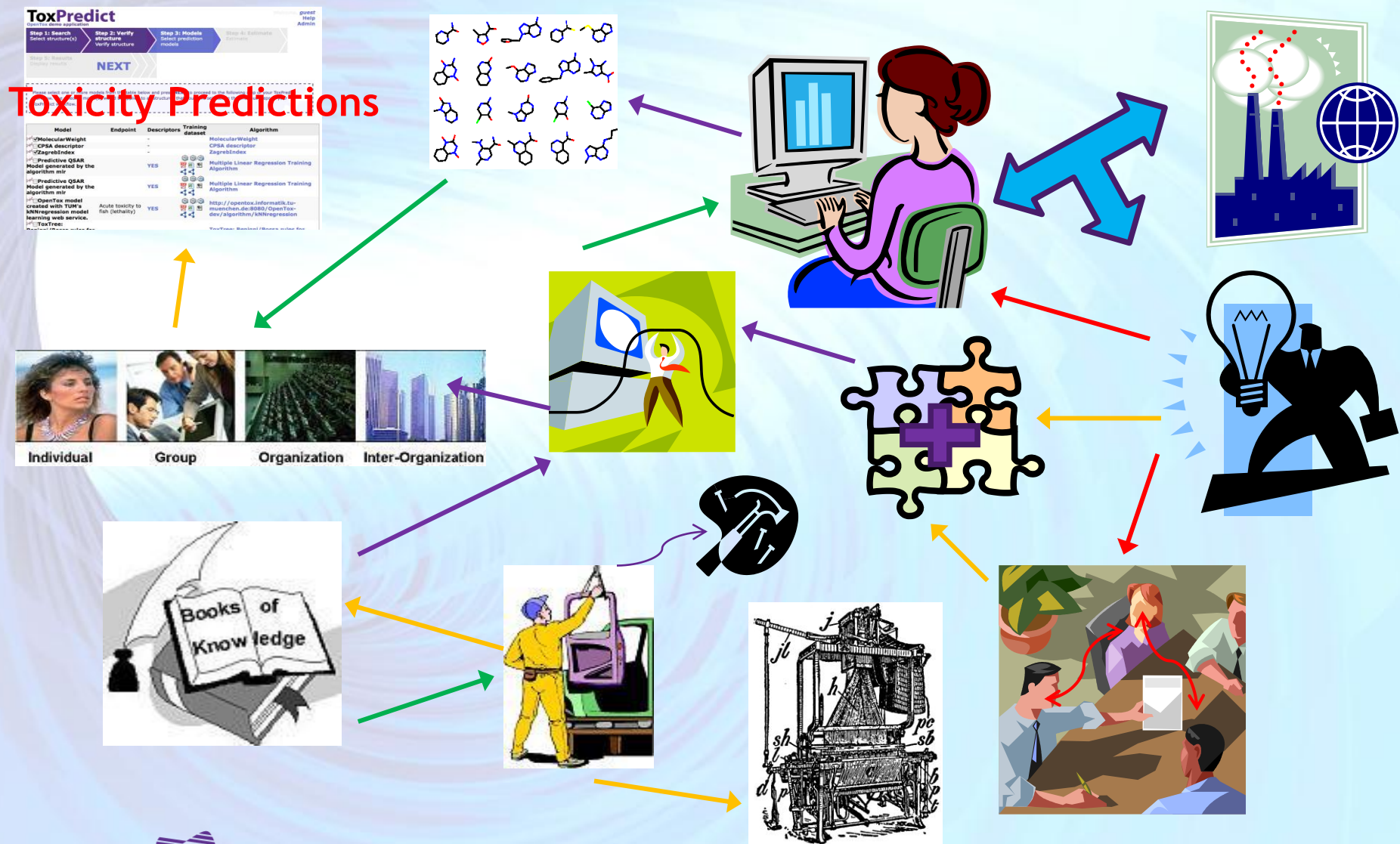
Solution created by Linked Open Data, Web Applications and Crowdsourcing



Haiti Earthquake Crisis Response (2010)

wiki.openstreetmap.org

Accelerating Knowledge Flows in Predictive Toxicology



OpenTox is an Integrating Framework

A diagram on the left side of the slide consists of three concentric purple semi-circles. The outermost semi-circle is the largest, the middle one is smaller, and the innermost one is the smallest. These semi-circles are positioned to the left of a table, with their right edges aligned with the table's columns. The table has three rows, each corresponding to one of the semi-circles. The first row is associated with the largest semi-circle, the second with the middle one, and the third with the smallest one. The table's first column contains the labels 'Framework', 'Diverse Access', and 'Interoperability' respectively. The second column contains bulleted lists of details for each category.

Framework

- Toxicity Data (Linked)
- *in silico* models
- Validation & Reporting
- Interpretation aids

Diverse Access

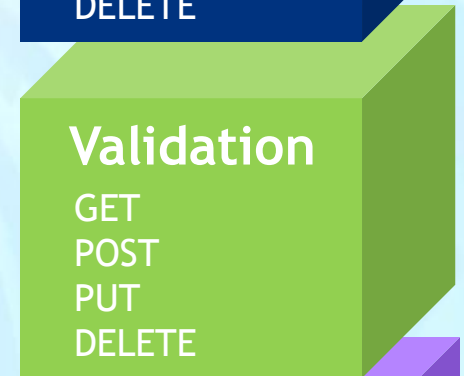
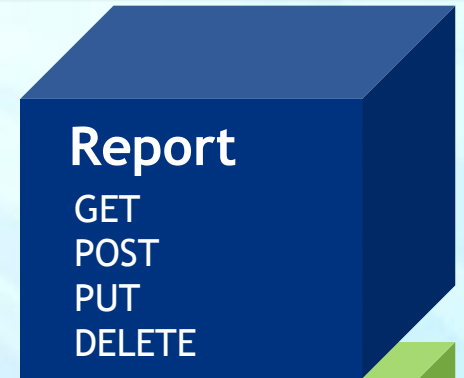
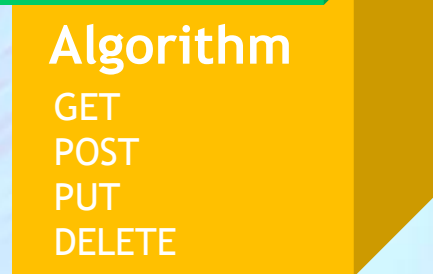
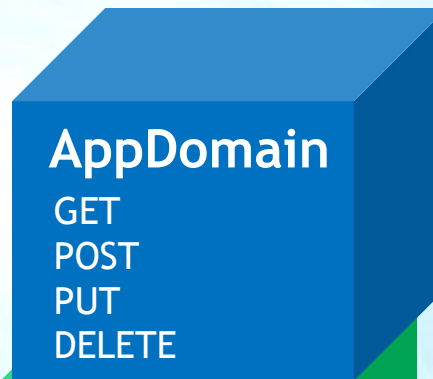
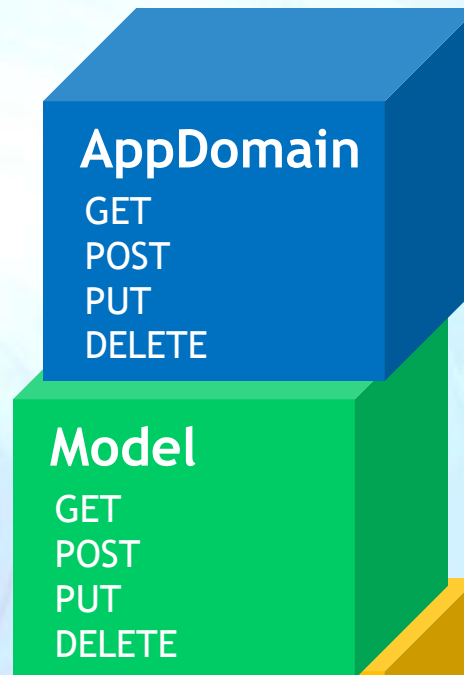
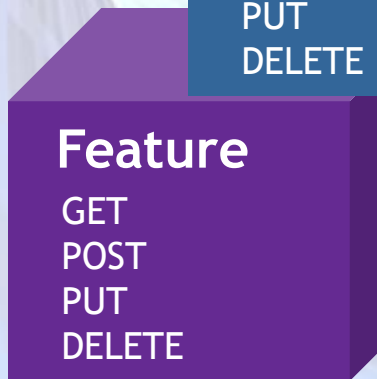
- Toxicologist, Biologist, Chemists
- Computational Scientists
- Interfaces for new algorithm development & integration

Interoperability

- Promote Standards
- Core Open Source Components
- Support Ontologies & Integration of Multiple Resources

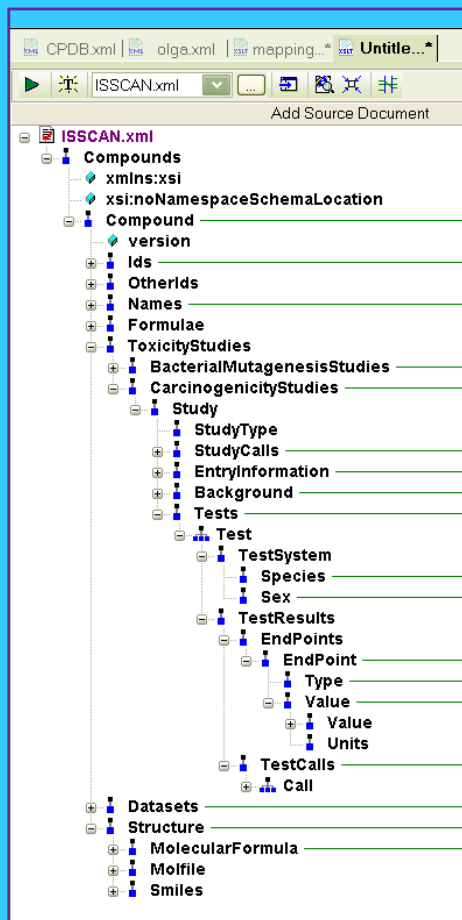
	OECD Principle	OpenTox addresses Validation Principles by...
1	Defined Endpoint	providing a unified source of well defined and documented toxicity data with a common vocabulary
2	Unambiguous Algorithm	providing transparent access to well documented models and algorithms as well as to the source code
3	Defined Applicability Domain	integrating tools for the determination of applicability domains during the validation of prediction models
4	Goodness-of-fit, robustness and predictivity	providing scientifically sound validation routines for the determination of errors and confidences
5	Mechanistic interpretation (if possible)	integrating tools for the inference, correlation or prediction of toxicological mechanisms and the recording of opinions and analysis in reports

Overview of Application Programming Interfaces



Toxicological Endpoint Ontology Development

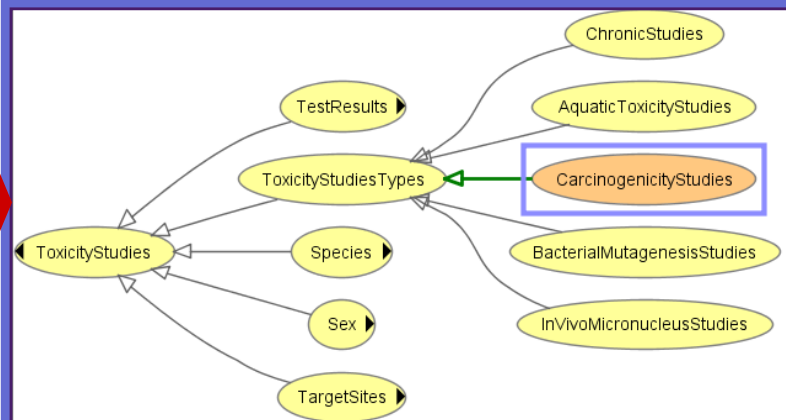
ToxML schema



Other publicly available resources:
DSSTox, GoReni (ITEM), ISSCAN ...

OpenTox Toxicological Endpoint Ontology

Ontology Development



Re-use of terms defined in
neighbouring ontologies (e.g. OBO)

Collaborative
Protégé
Environment

OpenToxipedia



Barry Hardy Log out Quicktools Site Setup Help

Site Map Accessibility Contact Data

Search Site

Home Toxicity Prediction OpenTox Blog People Partners Development OpenToxipedia
User Guidance Latest Entries A B C D E F G H I J K L M N O P Q R S T U V W
X Y Z by Categories Entries OpenToxipedia

You are here: Home » OpenToxipedia

Contents View Edit Rules Sharing History

Actions Display Add new... State: Published

OpenToxipedia

by Barry Hardy — last modified Sep 03, 2009 01:09 PM

OpenTox Community Resource for Toxicology Vocabulary and Ontology

OpenTox is supporting the creation and curation of OpenToxipedia, a community-based predictive toxicology knowledge resource. All members of the community are welcome to provide entries, suggested definition edits or additional information to entries in the resource.

OpenTox is supporting the application and development of the **ToxML** standard for representation of toxicology data, the **OECD principles for (Q)SAR model validation**, and the use of the **OECD HT** standard for regulatory reporting purposes.

OpenToxipedia provides here a Vocabulary Resource of toxicology terminology. We hope you find the resource useful and consider contributing to terms and their content.

Guidance for Vocabulary Resource entries



www.opentox.org/opentoxipedia



OpenTox: Databases

Chemical compounds - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://apps.ideaconsult.net:8180/ambit2/query/smarts?type=smiles&search=[*]OC(=O)[#6;H1]=[#6;H1]c1cccc1&t

Disable Cookies CSS Forms Images Information Miscellaneous Outline Resize Tools View Source Options

Chemical compounds

ToxPredict TTC Depiction Datasets Chemical compounds Similarity Substructure Algorithms References Features Templates Models Ontology RDF playground Help

ambit

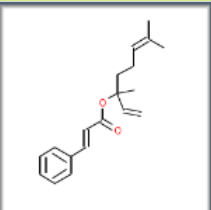
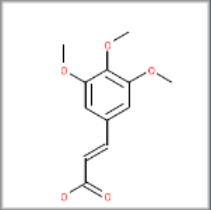
SMARTS

Keywords

Search for substructure and properties
This site and AMBIT REST services are under development!

Retrieve data

Search results SMARTS [*]OC(=O)[#6;H1]=[#6;H1]c1cccc1 Download as Max number of hits:

#	Compound	ECHA REGISTRATION DATE	ECHA CasRN	ECHA EC	ECHA Names	ECHA SYNON Names	ECHA SYNON Names	ECHA SYNON Names	ECHA SYNON Names	ECHA SYNON Names	ECHA SYNON Names
1		30.11.2010	78-37-5	201-110-3	linalyl cinnamate						
2		30.11.2010	90-50-6	201-999-8	3,4,5-trimethoxycinnamic acid						

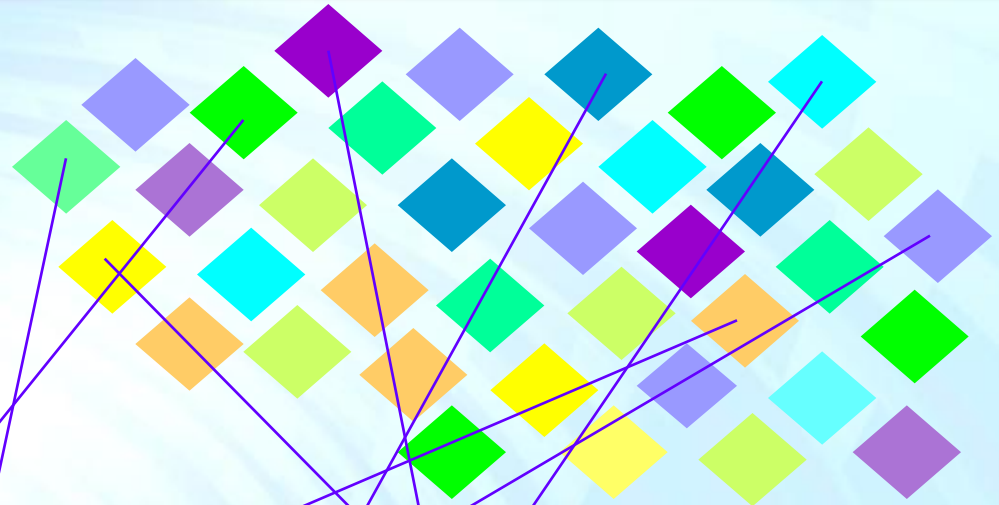
http://apps.ideaconsult.net

Creation of VO from Collaboration Pool

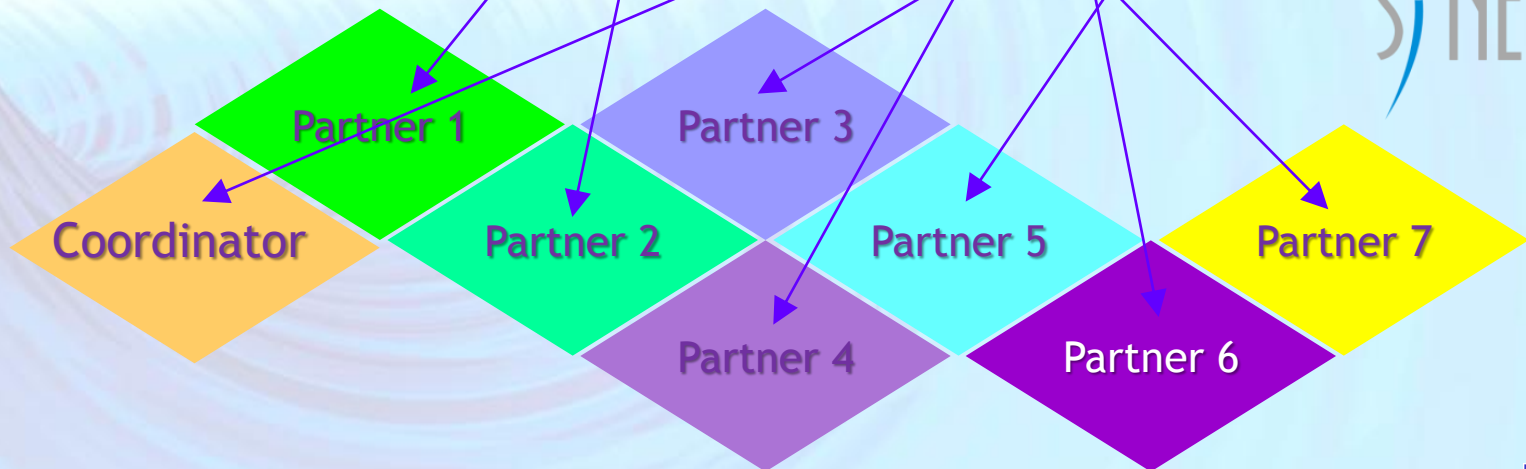
Network

Opportunity

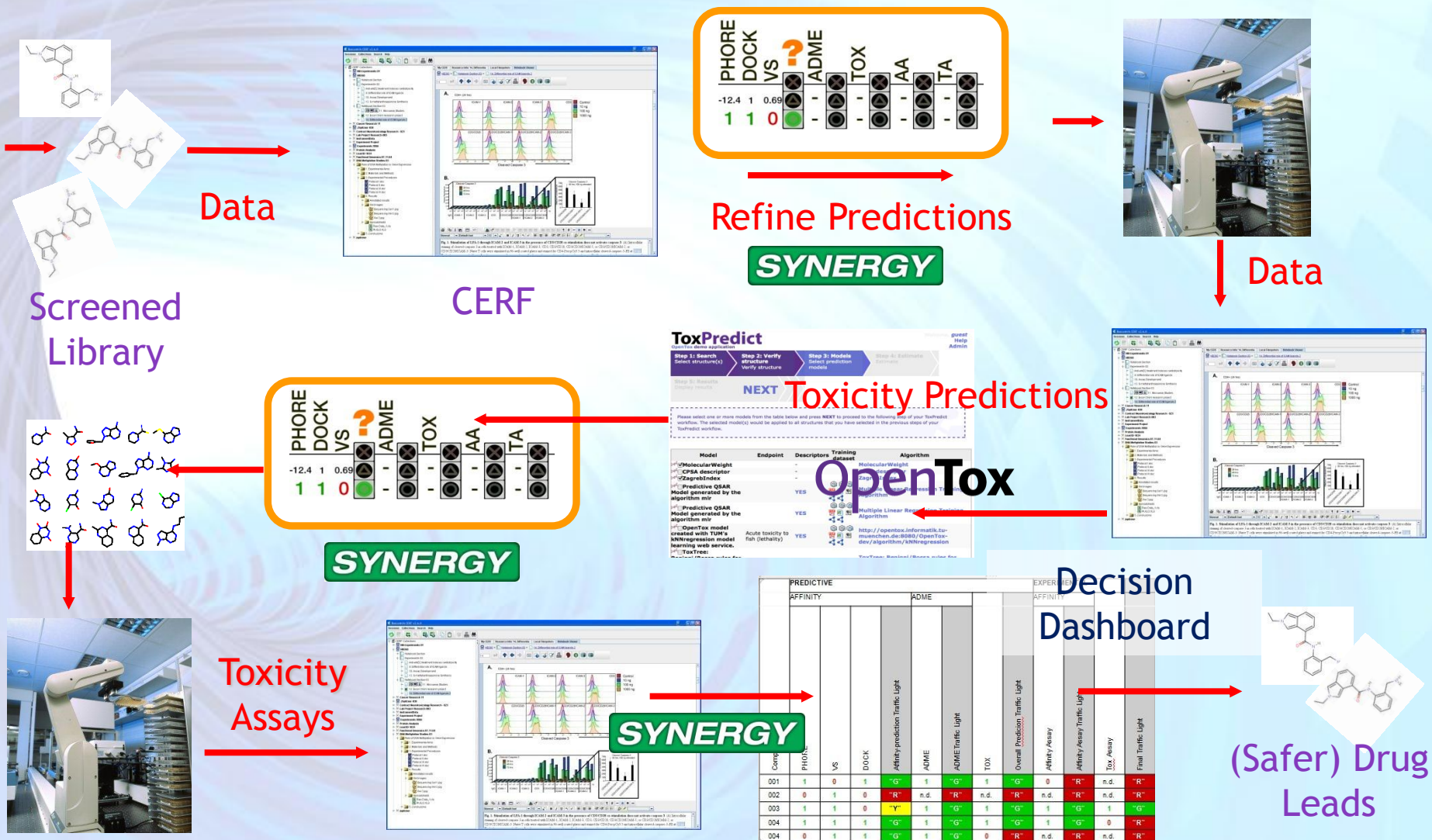
Call for Tender
Need for joint effort
Major project



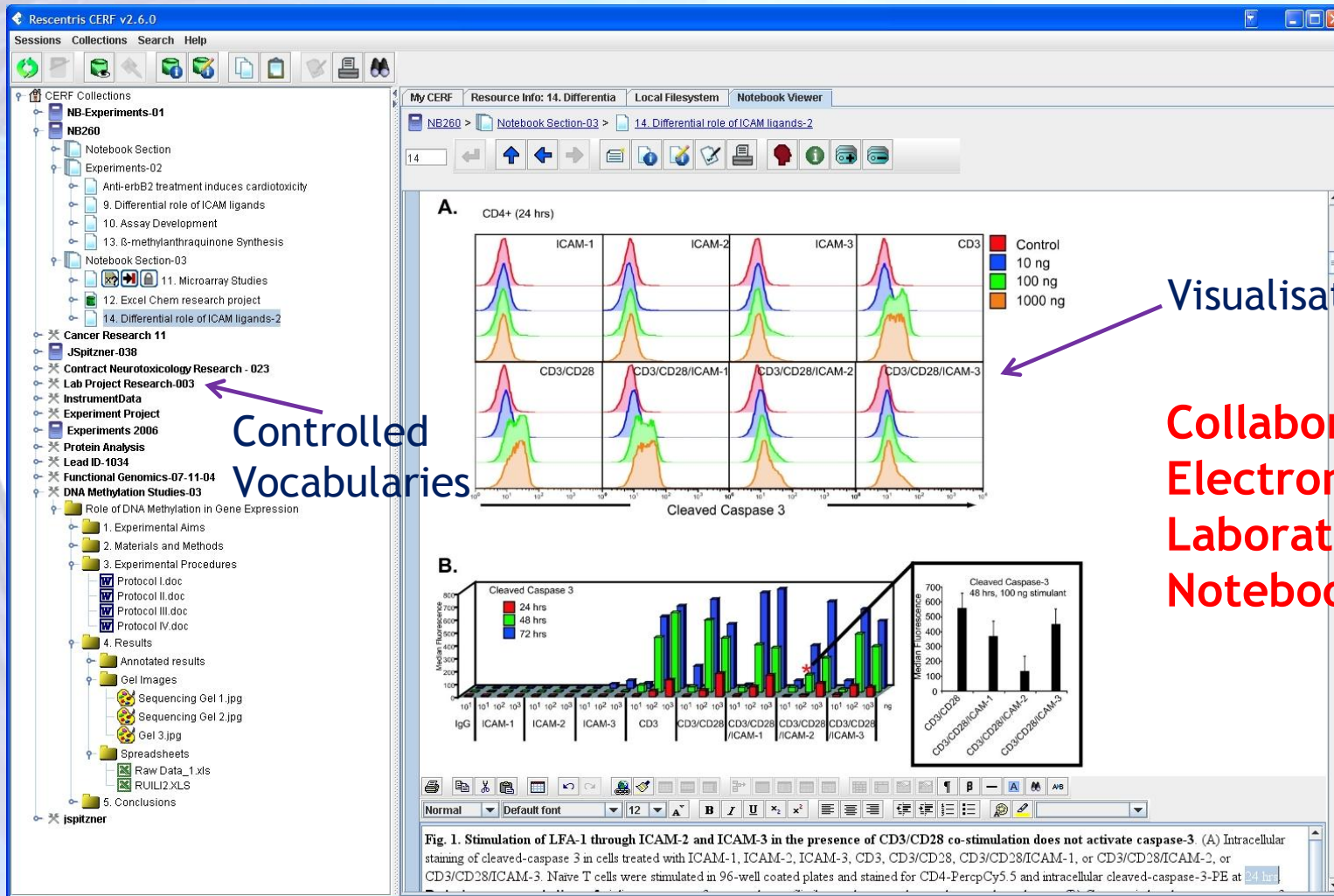
Virtual Organisation



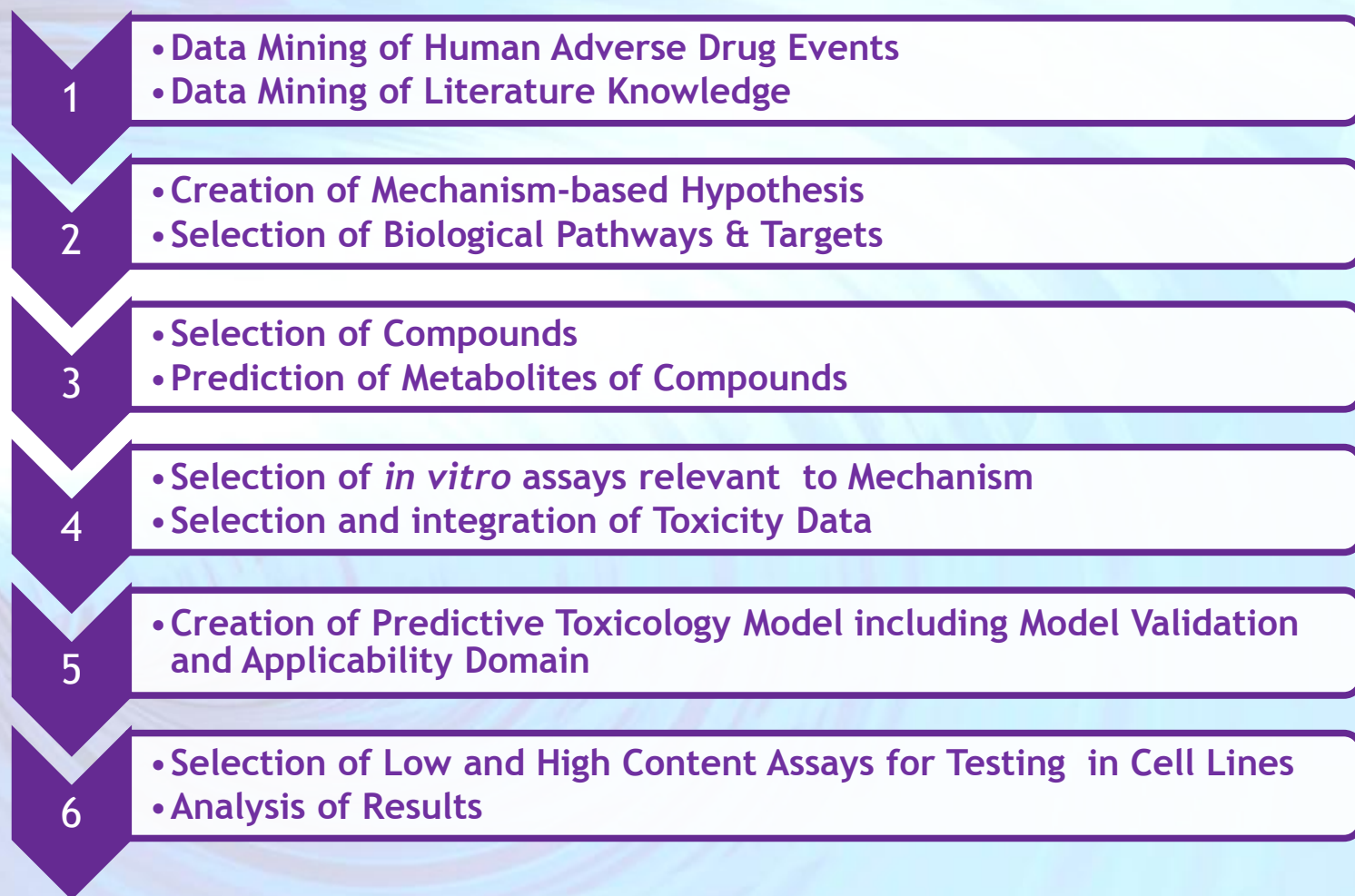
Synergy Drug Design Collaboration Pilot



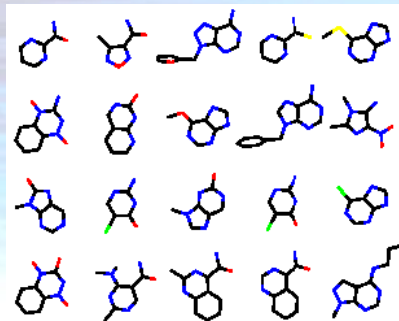
Recording of Collaborative R&D



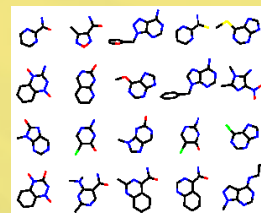
OpenTox - Synergy Predictive Toxicology VO Pilot Strategy Development & Case Study



1. A library of compounds is entered to the ELN



ELN

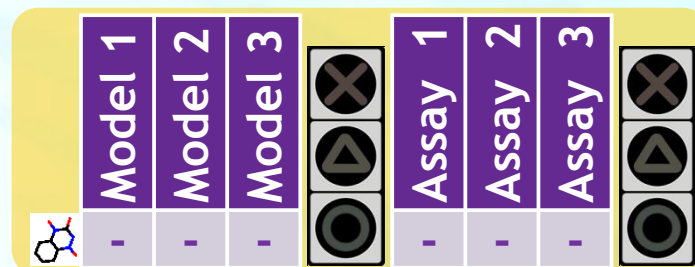


Synergy

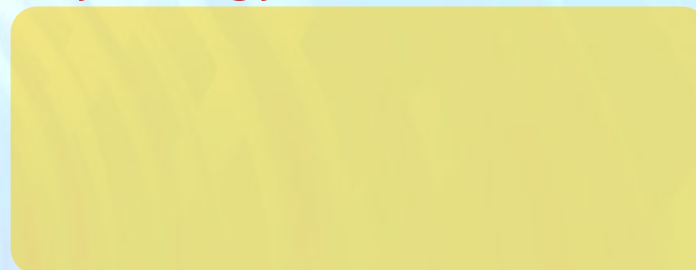
OpenTox

2. Each compound is assigned a data structure in ELN

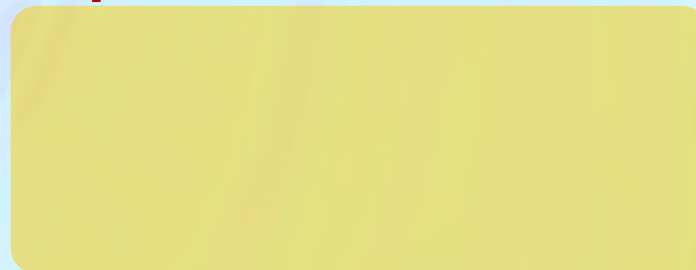
ELN



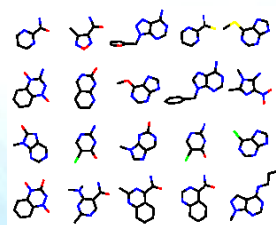
Synergy



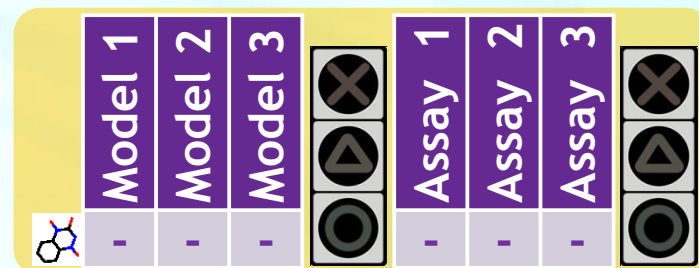
OpenTox



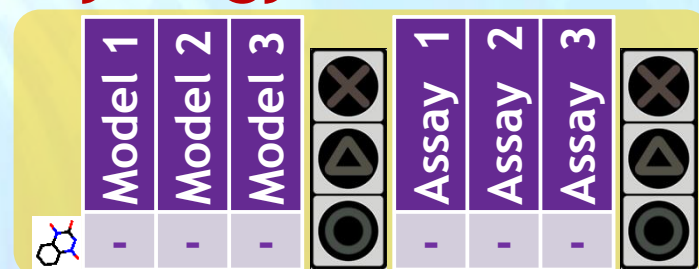
3. ELN passes compounds to OpenTox and SYNERGY



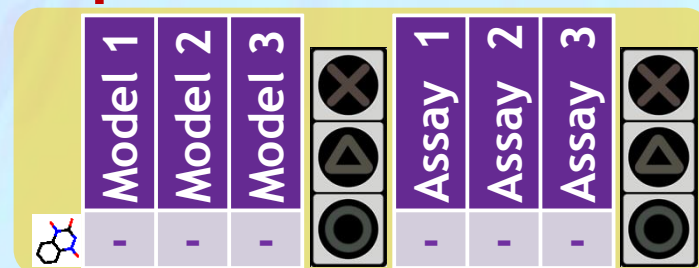
ELN



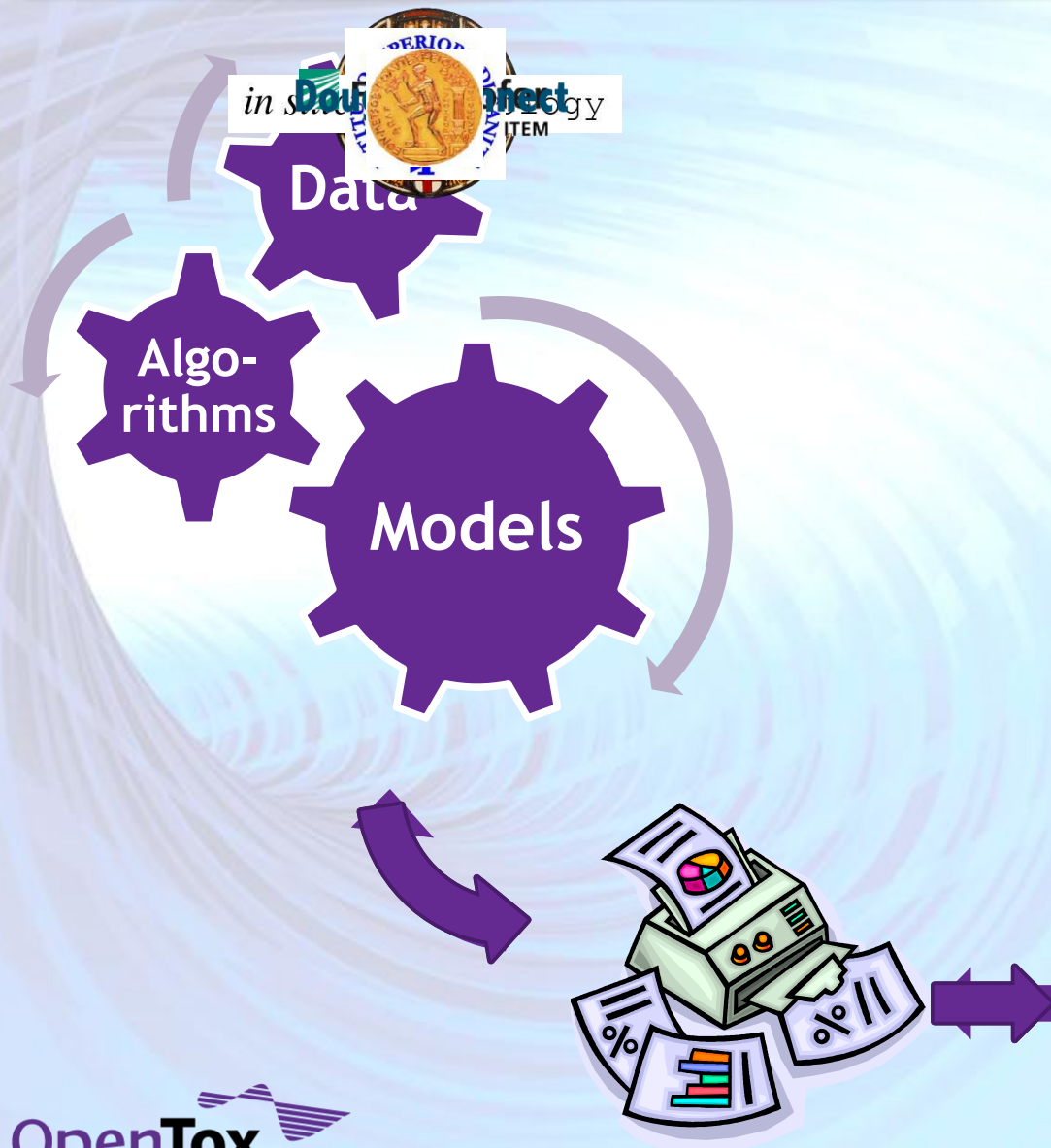
Synergy



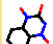






OpenTox



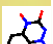






4. OpenTox computes toxicity predictions



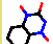






ELN

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	-	-	-		-	-	-	
	-	-	-		-	-	-	
	-	-	-		-	-	-	

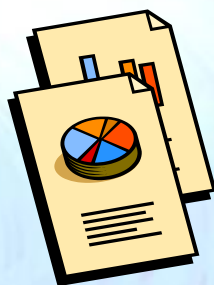
Synergy

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	-	-	-		-	-	-	
	-	-	-		-	-	-	
	-	-	-		-	-	-	

OpenTox

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	1	0	1		-	-	-	
	-	-	-		-	-	-	
	-	-	-		-	-	-	

5. OpenTox sends back a report to ELN



ELN

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	1	0	1		-	-	-	

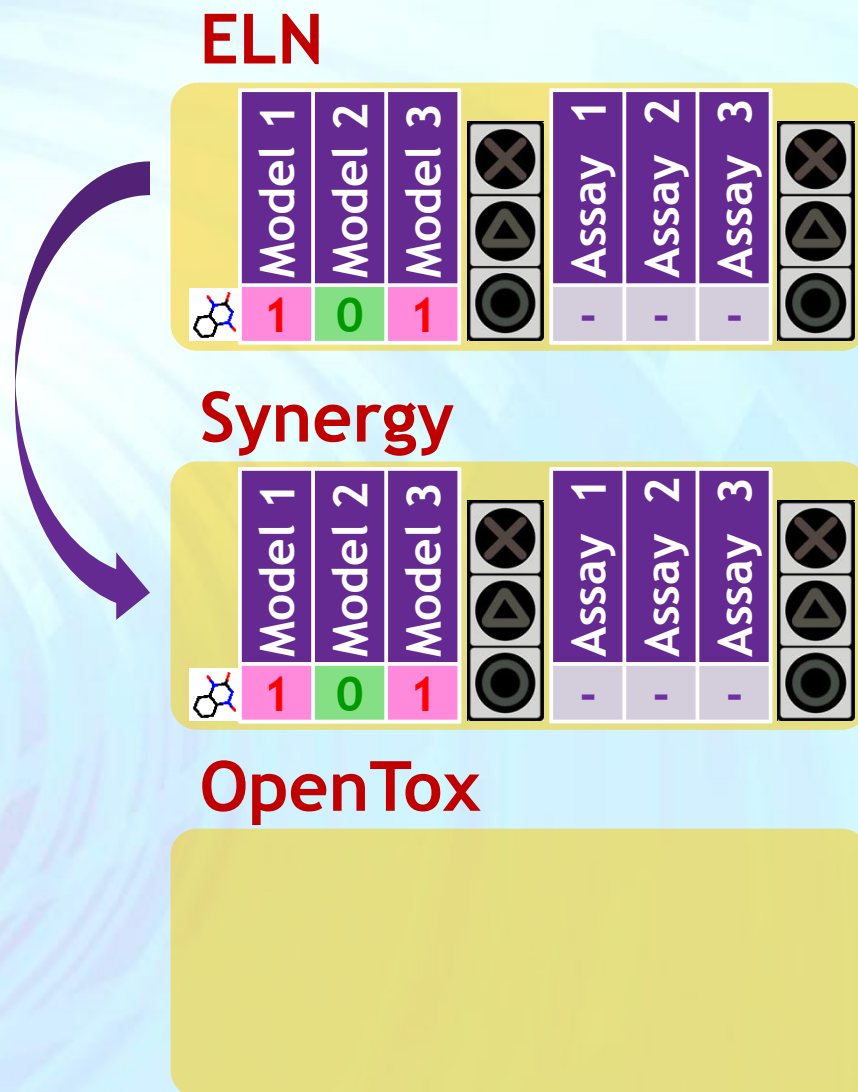
Synergy

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	-	-	-		-	-	-	

OpenTox

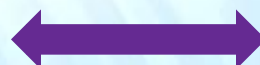
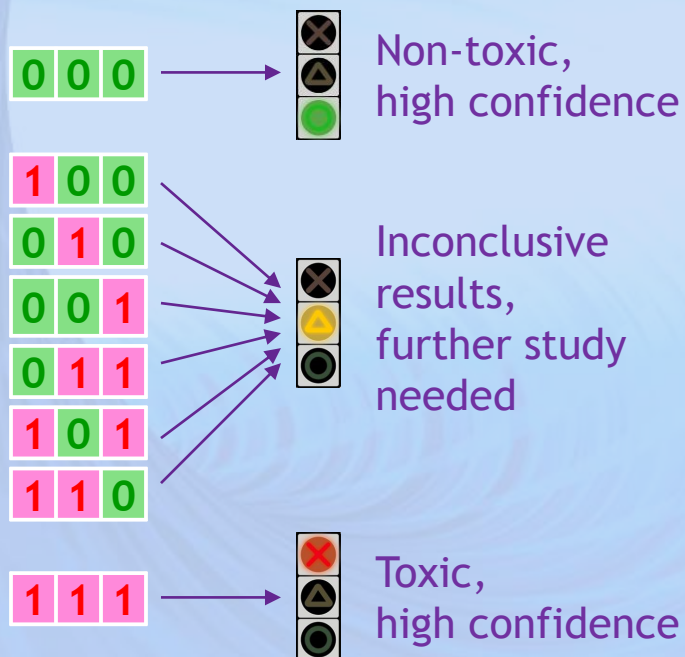
	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	1	0	1		-	-	-	

6. ELN sends the results to SYNERGY

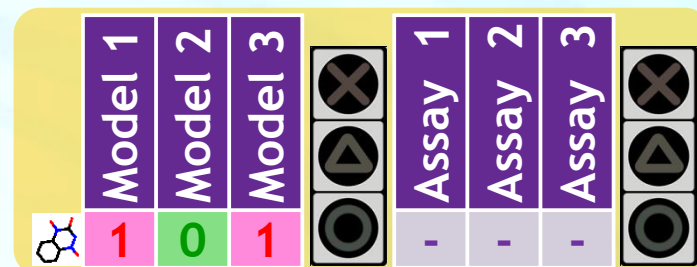


7. SYNERGY applies the Recommendation Rules

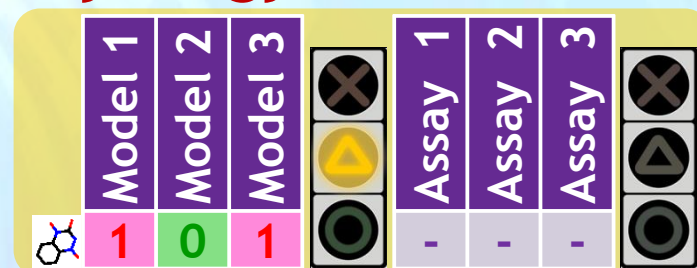
Recommendation Rules:



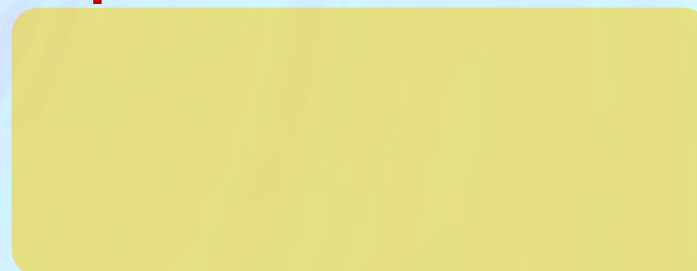
ELN



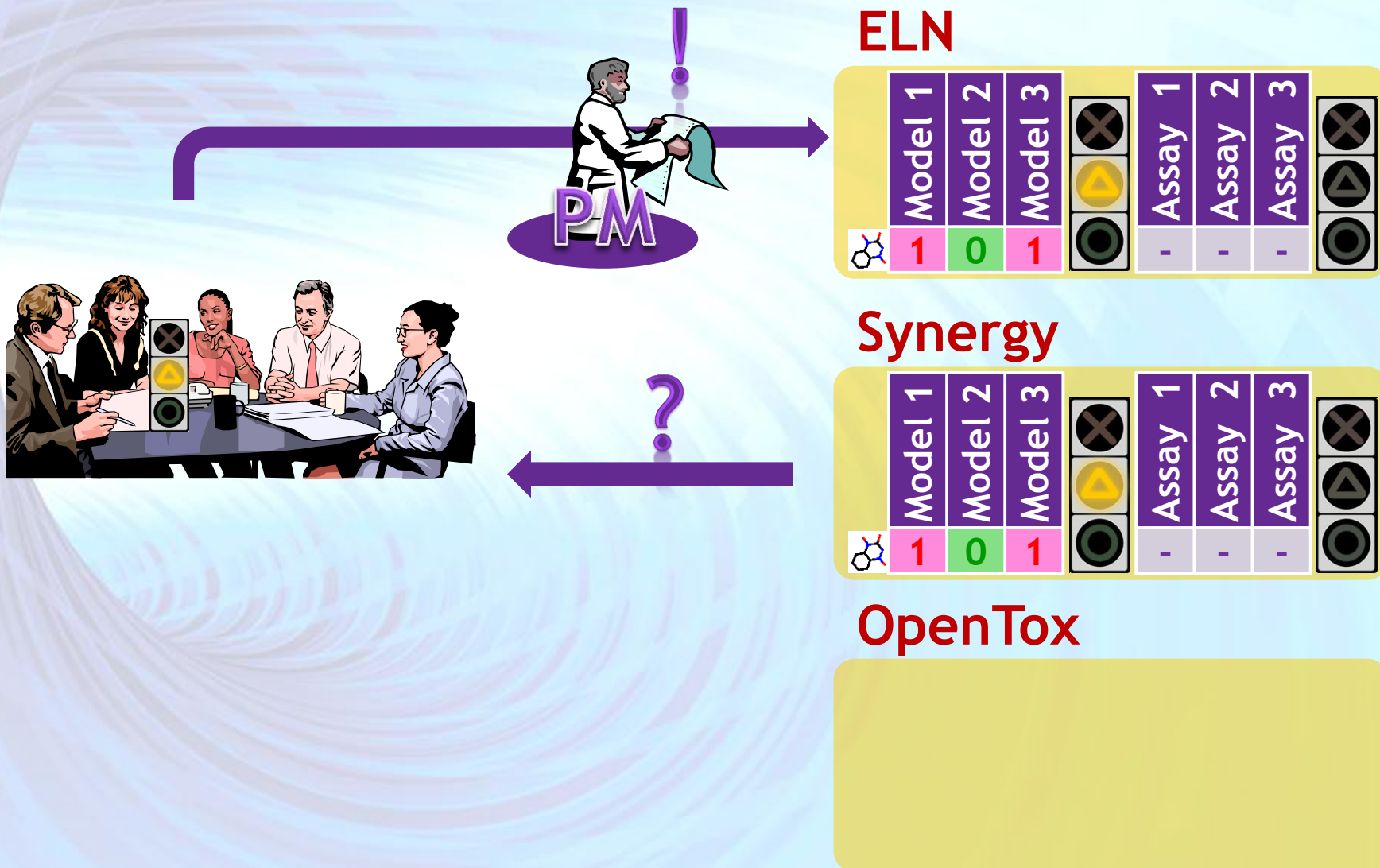
Synergy



OpenTox



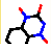









8. Inconclusive data → SYNERGY calls a meeting



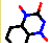






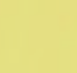

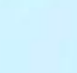
9. Experimental assays confirm toxicity



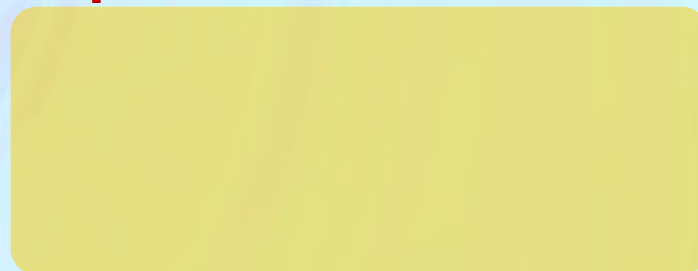
ELN

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	1	0	1		-	1	1	
								

Synergy

	Model 1	Model 2	Model 3		Assay 1	Assay 2	Assay 3	
	1	0	1		-	-	-	
								

OpenTox



Development and Use of Predictive Toxicology Applications

OpenTox Workshop
19 Sept. 2010, Rhodes, Greece

OpenTox Framework Design

Christoph Helma
(in silico toxicology)

Initial Motivation

- Predictive Toxicology applications need common components, e.g.
 - Access to datasets
 - Algorithms for descriptor calculation and model building
 - Validation routines
- These components have to be reimplemented 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 Authorisation: Protect confidential data

Requirements

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

Technological Choices

- Webservices
- 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

Interface Definitions

Description	Method	URI	Parameters	Result	Status codes
Retrieve SPARQL query results	GET	/ontology	? query =SPARQL_QUERY (mandatory)	RDF representation of the query results.	200,404,500
Predefined query to retrieve all models	GET	/ontology/models		RDF representation of all models.	
Predefined query to retrieve all endpoints	GET	/ontology/endpoints		RDF representation of all endpoints.	
Predefined query to retrieve all algorithms	GET	/ontology/algorithms		RDF representation of all algorithms.	
Submit SPARQL query and/or OpenTox service URL	POST	/ontology	uri []=URL of a OpenTox RDF resource query =SPARQL_QUERY	RDF representation of the query results, if query is specified. if uri [] is specified, the server retrieves a RDF representation and adds it to the RDF storage, thus making it available for the subsequent queries.	200,404,500,502

Ontologies

What?

- **Formal, shared conceptualization of a domain**

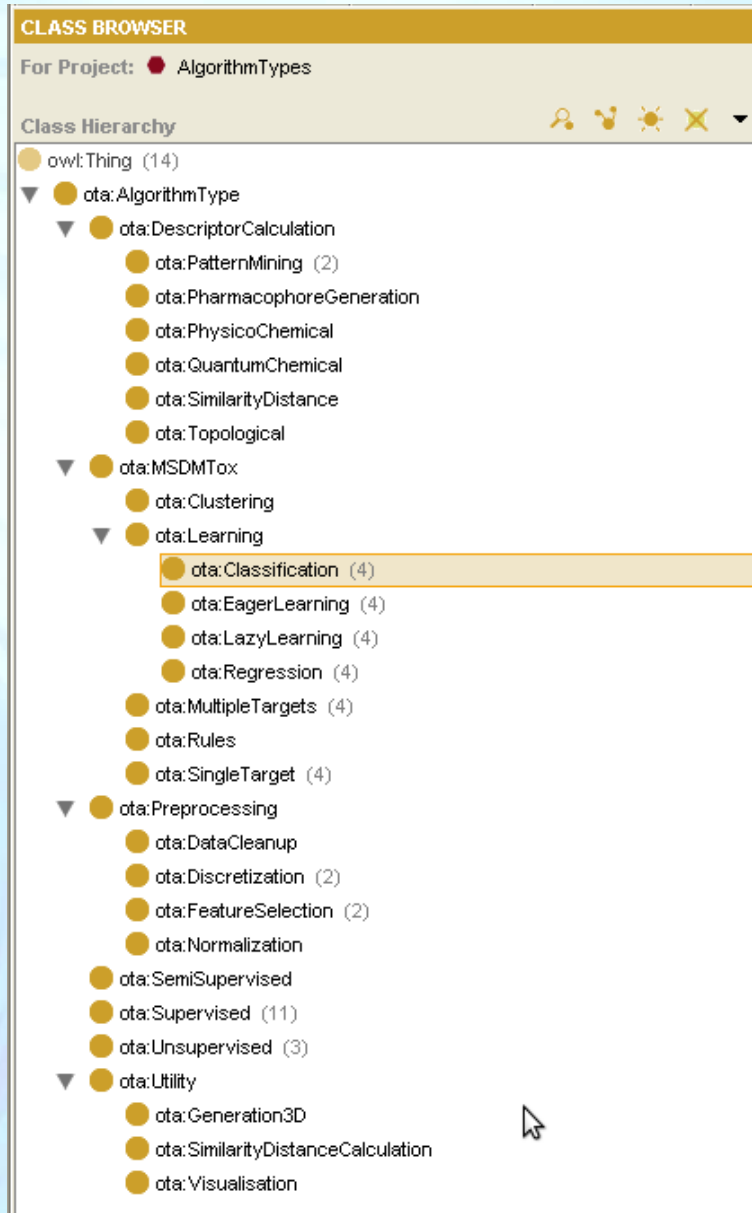
Why?

- 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

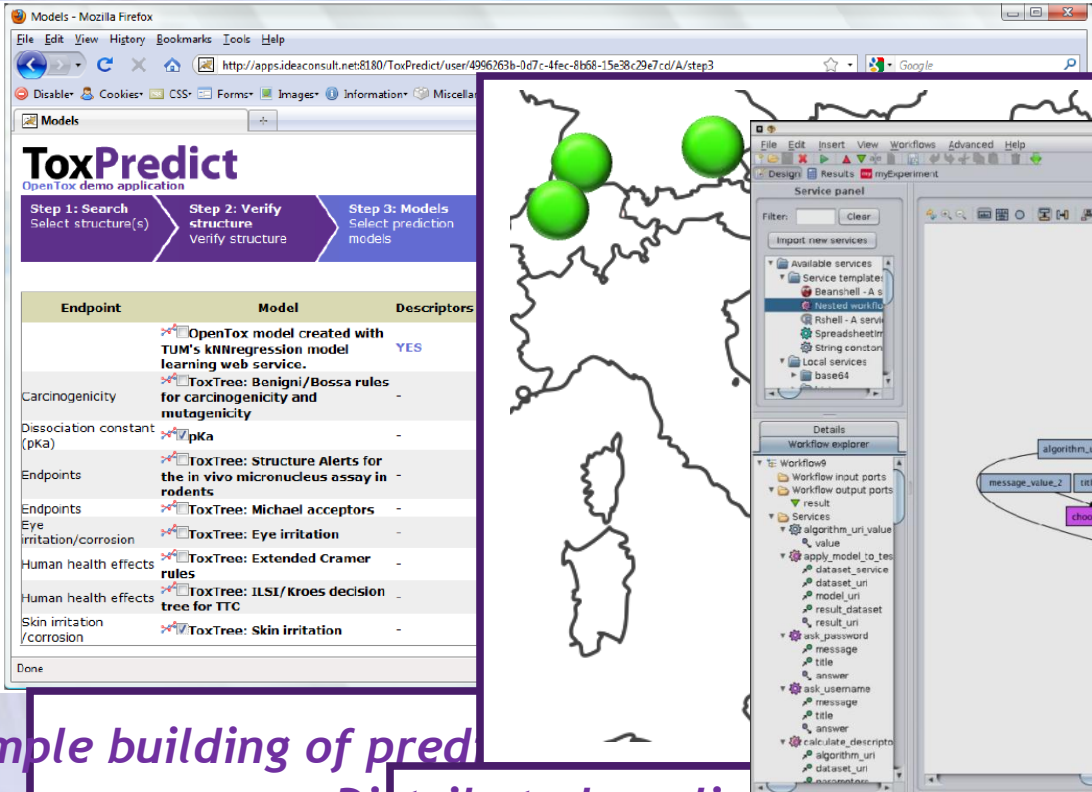
Ontologies

- Standards: **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



What can you do with OpenTox



ToxPredict
OpenTox demo application

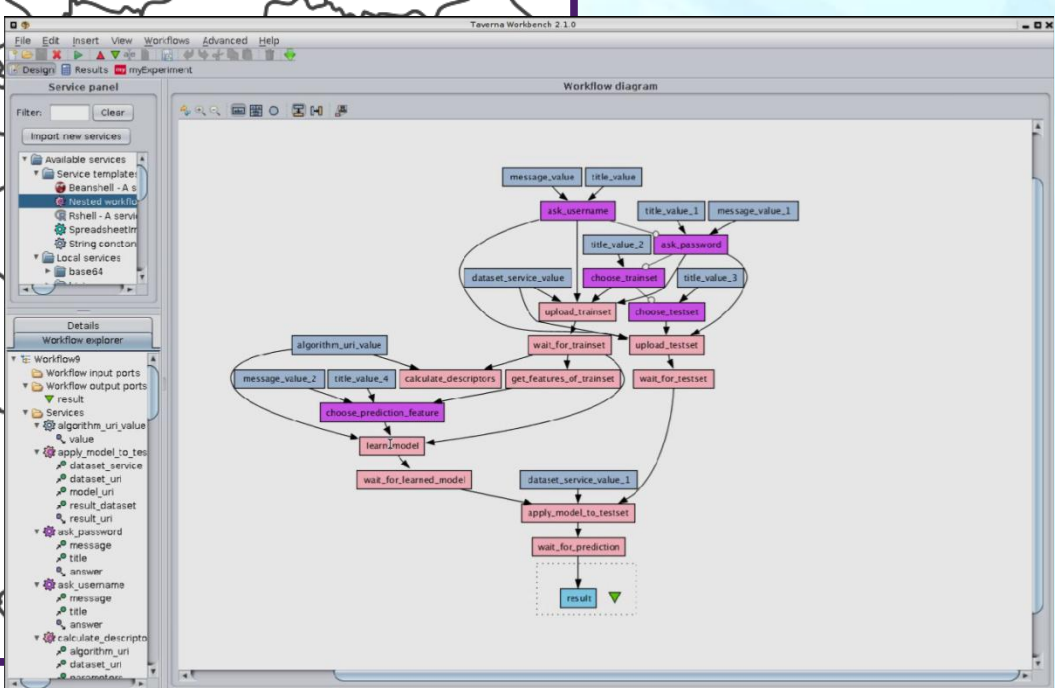
Step 1: Search
Select structure(s)

Step 2: Verify structure
Verify structure

Step 3: Models
Select prediction models

Endpoint	Model	Descriptors
Carcinogenicity	OpenTox model created with TUM's kNN regression model learning web service.	YES
Dissociation constant (pKa)	ToxTree: Benigni/Bossa rules for carcinogenicity and mutagenicity	-
Endpoints	ToxTree: Structure Alerts for the in vivo micronucleus assay in rodents	-
Endpoints	ToxTree: Michael acceptors	-
Eye irritation/corrosion	ToxTree: Eye irritation	-
Human health effects	ToxTree: Extended Cramer rules	-
Human health effects	ToxTree: ILSI/Kroes decision tree for TTC	-
Skin irritation/corrosion	ToxTree: Skin irritation	-

Done



Taverna Workflow editor

Workflow diagram

Service panel

Workflow explorer

Workflow diagram details:

- Input: message_value, title_value, title_value_1, message_value_1
- Process: ask_username, title_value_2, ask_password, title_value_3
- Process: dataset_service_value, choose_trainset, choose_testset
- Process: upload_trainset, upload_testset
- Process: wait_for_trainset, wait_for_testset
- Process: calculate_descriptors, get_features_of_trainset
- Process: choose_prediction_feature, learn_model
- Process: wait_for_learned_model
- Process: dataset_service_value_1, apply_model_to_testset
- Process: wait_for_prediction
- Output: result

Simple building of prediction applications based on a wide range of data methods and data integration into workflow systems for computational biology

Summary

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

Development and Use of Predictive Toxicology Applications

An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Application Programming Interfaces

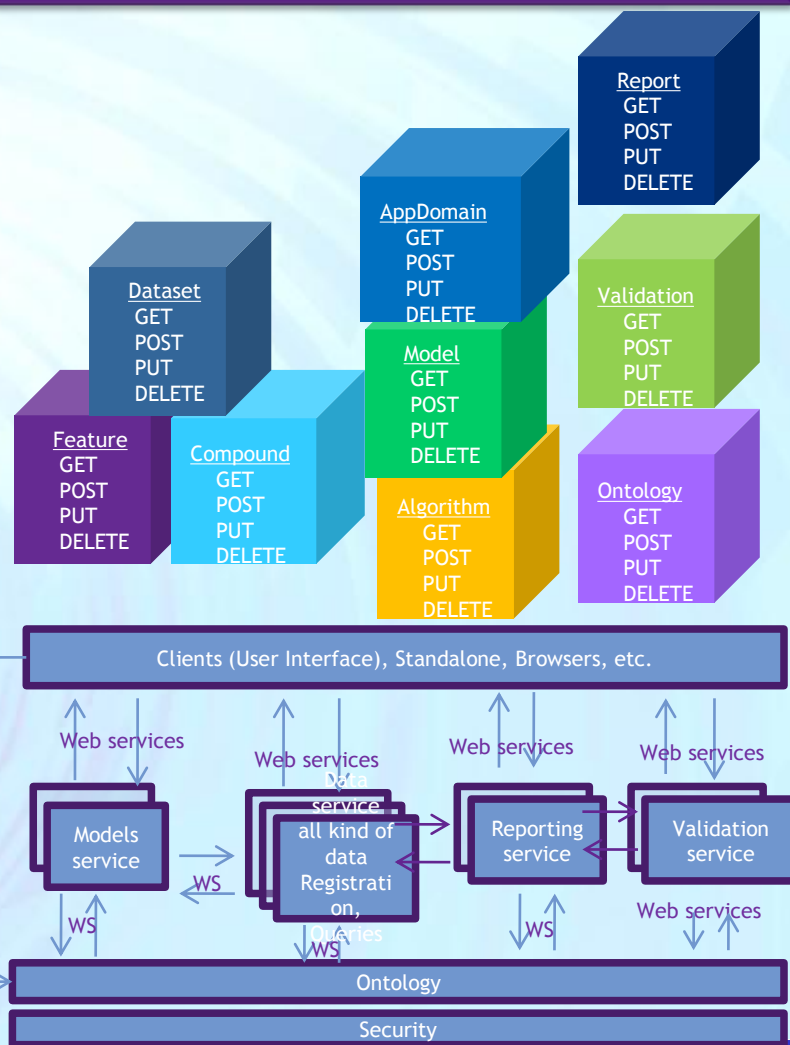
presented by Nina Jeliazkova
(Ideaconsult Ltd., Bulgaria)

Framework design rationales

User Requirements		Software Requirements
Unambiguous data	⇒	<i>formal way of representing information about data</i>
Unambiguous access	⇒	<i>well-defined interfaces</i>
Transparency of computational tools	⇒	<i>formal way of representing information about methods, well-defined interfaces</i>
Variety of user groups	⇒	<i>simplicity and modularity of design</i>
Need to integrate various resources (e.g., databases, prediction methods, models, ...) to make meaningful predictions	⇒	<i>distributed architecture, interoperability</i>
Need to integrate biological information	⇒	<i>again, modularity of design, extensibility</i>

The framework

- OpenTox API
 - The way applications talk to each other
 - The way developers talk to applications
 - <http://opentox.org/dev/apis/api-1.1>
- The basic building blocks:
 - data, chemical structures, algorithms and models.
- Functionality offered
 - build models,
 - apply models,
 - validate models,
 - access and query data in various ways.
- Technologies
 - REST style web services
 - RDF for description of resources
 - Links to existing and newly developed ontologies (mainly to describe metadata) about resources



Representational State Transfer (REST)

A software architecture style, defined by Roy Fielding in his [PhD thesis \(2000\)](#). Many services worldwide offer REST API. There are (currently) no standards for RESTful applications, but merely design guides.

Design principles:

- Resource oriented
 - Every object (resource) is named and addressable (e.g. HTTP URL) Example: <http://example.opentox.com/model/myBestModel> , <http://example.opentox.com/compound/50-00-0>
 - RESTfull API design starts by identifying most important objects and groups of objects, supported by the software system and proceeds by defining URL patterns.
- Transport protocol
 - HTTP is the most popular choice of transport protocol, but other protocols can be used as well
- Operations
 - All resources (nouns) support the same fixed and universal number of operations (verbs). HTTP (GET, POST, PUT, DELETE) operations are the common choice, when the transport protocol is HTTP.
- Hypermedia as the Engine of Application State
 - All resources should be reachable via a single (or minimum) number of entry points into RESTful applications. Thus, a representation of a resource should return hypermedia links to related resources
- Error codes (for each resource/operation pair)
 - HTTP status codes (e.g. 200 OK, 400 Bad Request, 404 Not found, etc.) are usually used

OpenTox resources (1)

OpenTox considers the following set of entities as essential building blocks:

- Structures of **chemical compounds**
- **Properties and identifiers** of chemical compounds
- **Datasets** of chemical compounds and various properties (measured or calculated)
- **Algorithms**
 - Data processing algorithms
 - Algorithms generating certain values, based on chemical structure (e.g. descriptor calculation)
 - Data preprocessing (e.g. Principal component analysis, feature selection)
 - Structure processing (e.g. structure optimization)
 - Algorithms, relating set of structures to another set of structures (e.g. similarity search or metabolite generation)
 - Machine learning algorithms
 - Supervised (e.g. Regression, Classification)
 - Unsupervised (e.g. Clustering)
 - Prediction algorithms, defined by experts (e.g. series of structural alerts, defined by human experts , not derived by learning algorithms)

OpenTox resources (2)

- **Models** are generated by respective algorithms, given specific parameters
 - 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.
- **Validation** provides procedures independent of model building facilities (e.g. crossvalidation) and generates relevant statistics.
- **Reports**
 - Various types of reports might be generated, using building blocks above (e.g. validation report can be generated using validation object, a model and a dataset).
- In addition, the following components are introduced:
 - **Task** (asynchronous processing of computationally intensive tasks)
 - **Authentication and authorization** (Ensuring secure access to sensitive resources)
 - **Ontology service** (provides an RDF storage and SPARQL endpoint for resources registration)

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/{modelid}
Algorithm	OpenTox algorithm services	http://host:port/algorithm/{algorithmid}
Validation, Report	A validation corresponds to the validation of a model on a test dataset.	http://host:port/validation/{validationid} 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 REST operations

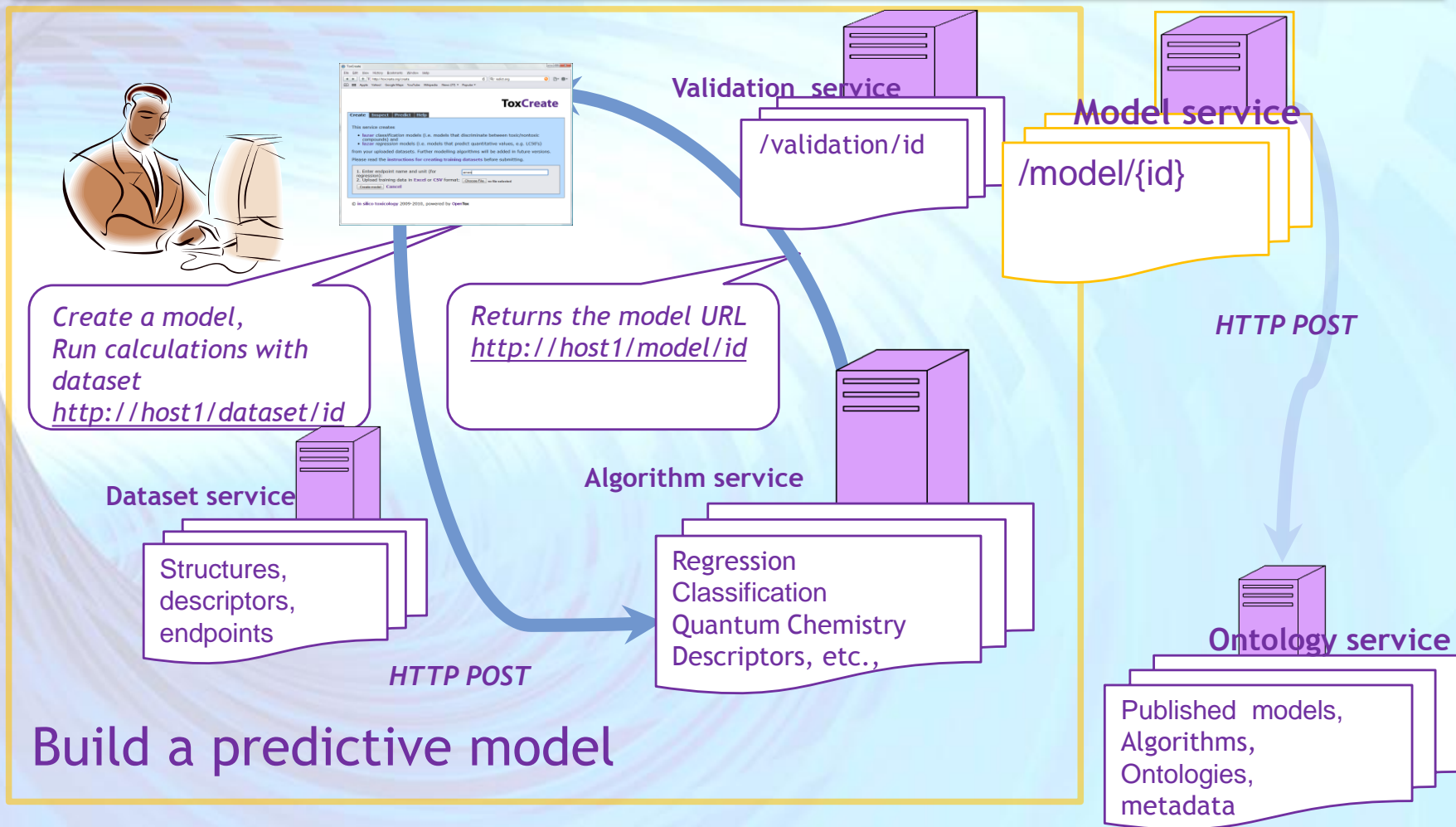
Individual resources (e.g. a dataset or a model)

- URI template <http://host:port/{resource}/{resourceid}> , e.g. http://host:port/model/{model_id} or http://host:port/dataset/{dataset_id}
- GET - retrieve representation of the resource
- PUT - update representation of the resource
- POST :
 - replace representation of the resource with a new one (e.g. replace the dataset with new content)
 - initiate calculations, based on this resource (e.g. submit dataset URI to an algorithm resource and obtain a model URI as a result)
- DELETE - delete the resource

Collections of resources (e.g. list of all available models, or datasets)

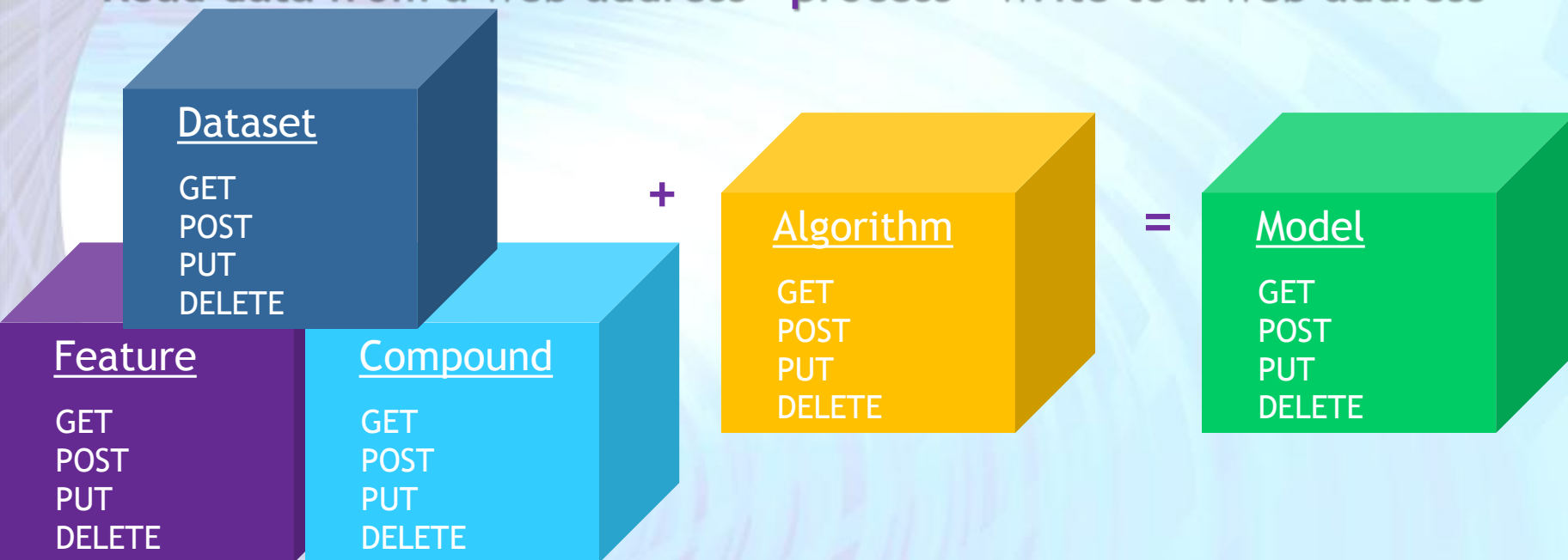
- URI template <http://host:port/{resource}> , (e.g. <http://host:port/model> or <http://host:port/dataset>)
- GET - retrieve representation of multiple resources (e.g. retrieve all available algorithms)
- PUT - N/A
- POST - create new resource and return its URI (e.g. create a new dataset by submitting new dataset content to the dataset service)
- DELETE - N/A

Build a predictive model



Uniform approach to models creation

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



<http://myhost.com/algorithm/neuralnetwork>

<http://myhost.com/dataset/trainingset1>

<http://myhost.com/model/predictivemodel1>

Use an algorithm to build a model

- An algorithm is applied by submitting HTTP POST to the algorithm URI and providing required parameters.
- A common required parameter is **dataset_uri=http://host:port/dataset/{datasetid}**, which specifies the data set to be operated on.
- HTTP POST in REST style services returns URI of the result, and not the content of the result.
- The algorithm services are designed to store the results into a dataset service and return the URL of the resulted dataset.
- In case of slow calculations a Task URI, instead of the dataset URI is returned

```
$ curl -H "Accept:text/uri-list" -X POST -d
'dataset_uri=http://apps.ideaconsult.net:8080/ambit2/dataset/1037' -d
'prediction_feature=http://apps.ideaconsult.net:8080/ambit2/feature/26
701' -d
'dataset_service=http://apps.ideaconsult.net:8080/ambit2/dataset'
http://opentox.informatik.tu-muenchen.de:8080/OpenTox-
dev/algorithm/J48 -iv
* Connected to opentox.informatik.tu-muenchen.de (131.159.28.16) port
8080 (#0)
> POST /OpenTox-dev/algorithm/J48 HTTP/1.1
>> Host: opentox.informatik.tu-muenchen.de:8080
> Accept: */*
> Content-Type: application/x-www-form-urlencoded
< HTTP/1.1 202 Accepted
< Date: Sat, 31 Jul 2010 14:46:38 GMT
< Location: http://opentox.informatik.tu-muenchen.de:8080/OpenTox-
dev/task/acdf6eac-d5a2-402c-a4e2-06cd7e3ca1b5
< Accept-Ranges: bytes
< Server: Noelios-Restlet-Engine/1.1.snapshot
< Content-Type: text/uri-list; charset=ISO-8859-1
< Content-Length: 99
<
* Connection #0 to host opentox.informatik.tu-muenchen.de left intact
* Closing connection #0
http://opentox.informatik.tu-muenchen.de:8080/OpenTox-
dev/task/acdf6eac-d5a2-402c-a4e2-06cd7e3ca1b5
```

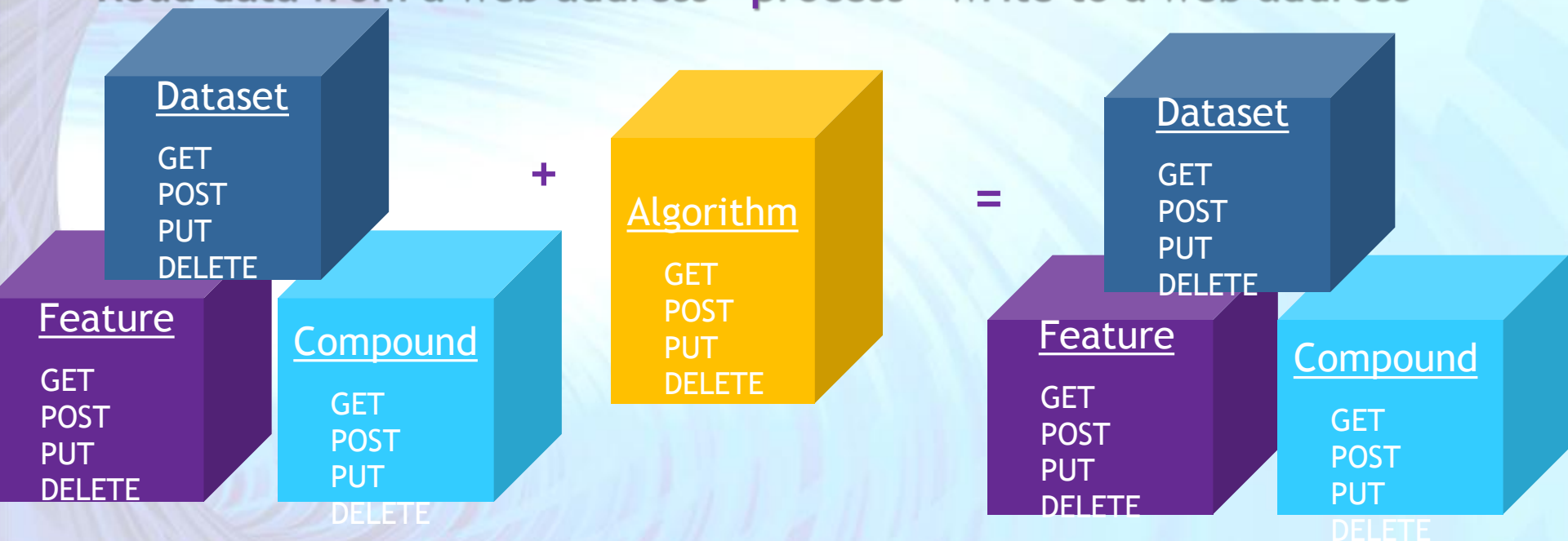
Resources: The model

- When task URI is returned, the returned status code is HTTP 202 Accepted, instead of HTTP 200 OK.
- This tells the client the processing is not completed and the client needs to poll the task URI until OK code is returned
- The final result, returned by Example 25 is the URI of the new model http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/model/TUMOpenToxModel_j48_48.
- To obtain prediction results POST a dataset to the model URI

```
$ curl -iv -H "Accept:text/uri-list" http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/task/acdf6eac-d5a2-402c
* About to connect() to opentox.informatik.tu-muenchen.de port 8080 (#0)
* Trying 131.159.28.16... connected
* Connected to opentox.informatik.tu-muenchen.de (131.159.28.16) port 8080 (#0)
> GET /OpenTox-dev/task/acdf6eac-d5a2-402c-a4e2-06cd7e3ca1b5 HTTP/1.1
> User-Agent: curl/7.18.2 (x86_64-pc-linux-gnu) libcurl/7.18.2 OpenSSL/0.9.8g
zlib/1.2.3.3 libidn/1.8 libssh2/0.18
> Host: opentox.informatik.tu-muenchen.de:8080
> Accept:text/uri-list
>
< HTTP/1.1 200 OK
< Date: Sat, 31 Jul 2010 14:47:22 GMT
Date: Sat, 31 Jul 2010 14:47:22 GMT
< Location: http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/model/TUMOpenToxModel\_j48\_48
< Vary: Accept-Charset, Accept-Encoding, Accept-Language, Accept
< Accept-Ranges: bytes
< Server: Noelios-Restlet-Engine/1.1.snapshot
< Content-Type: text/uri-list; charset=ISO-8859-1
< Content-Length: 86
<
* Connection #0 to host opentox.informatik.tu-muenchen.de left intact
* Closing connection #0
http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/model/TUMOpenToxModel\_j48\_48
```

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

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



<http://myhost.com/algorithm/{descriptorX}>

<http://myhost.com/dataset/trainingset1>

<http://myhost.com/dataset/results>

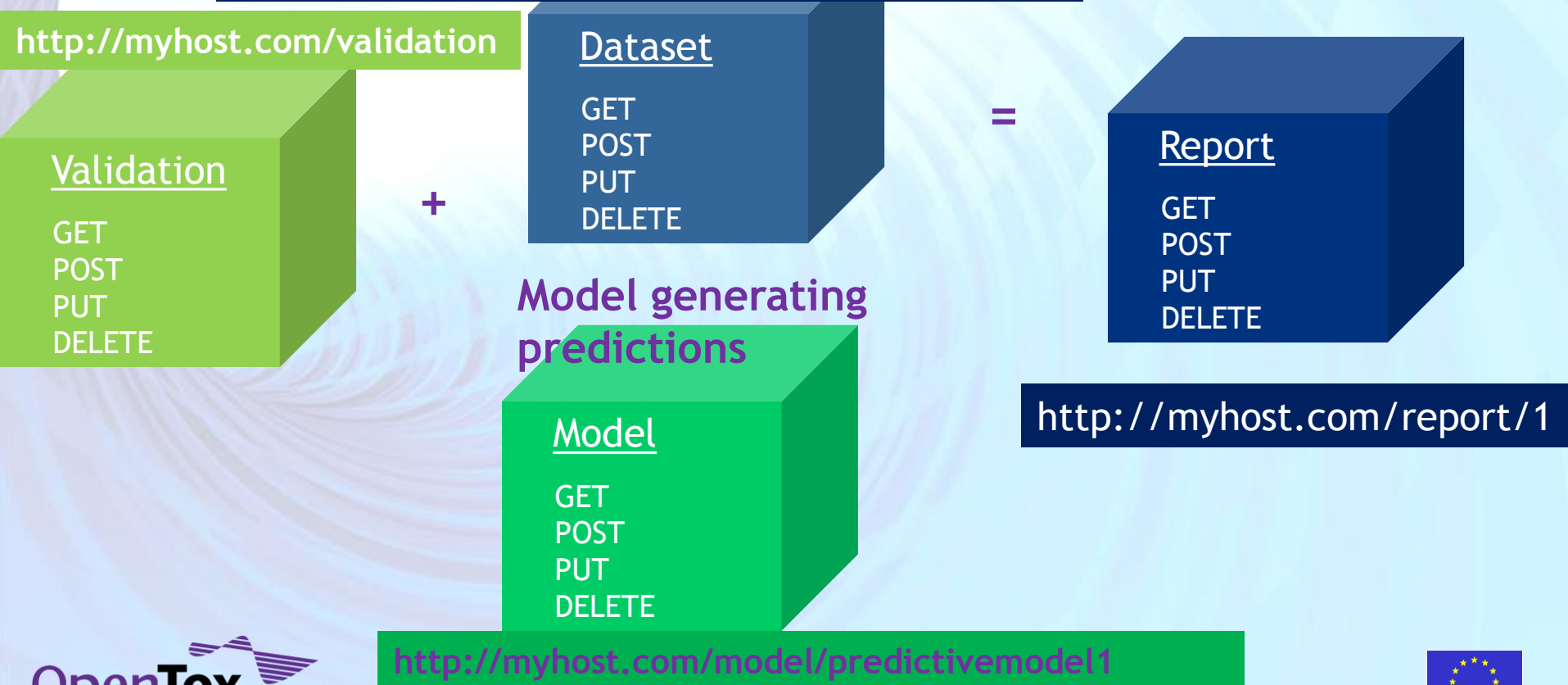
Uniform approach to models validation and report generation

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

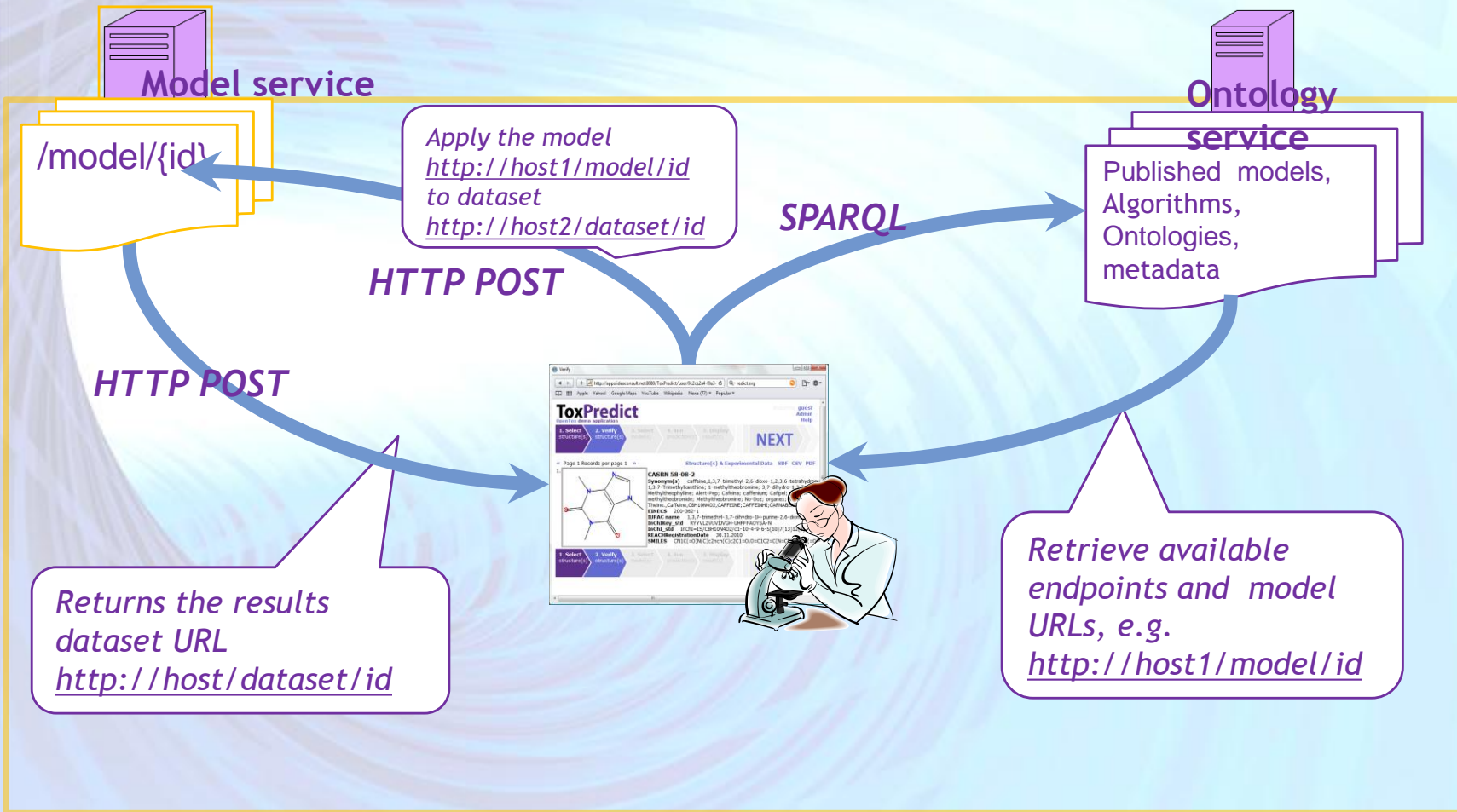
<http://myhost.com/dataset/trainingset1>

<http://myhost.com/dataset/predictedresults1>

Validation report

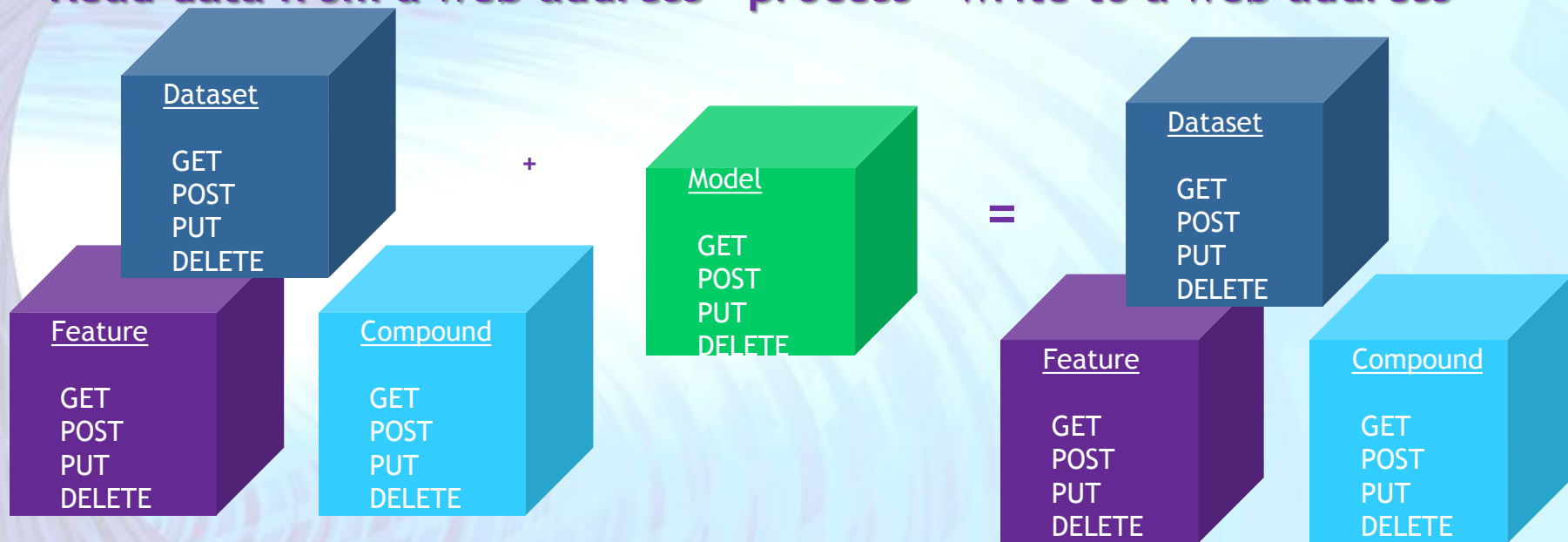


Apply predictive models



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>

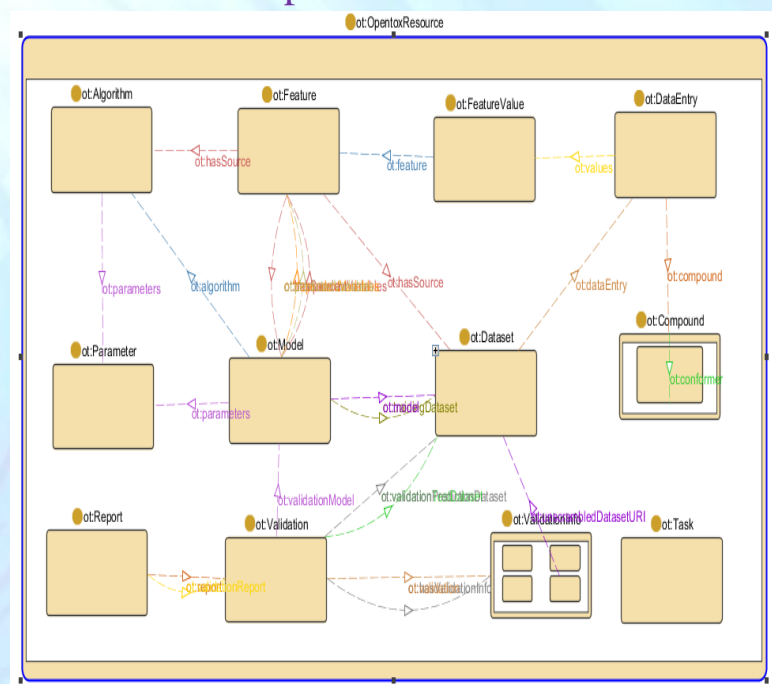
RDF - Resources representation

- The [opentox.owl](http://opentox.org/api/1.1/opentox.owl) ontology
 - A common OWL data model of all OpenTox resources
 - Describes OpenTox resources
 - Describes relationships between them
 - Generates object's RDF representations.
- RDF/XML representation is mandatory for OpenTox resources.
- Uniform approach to data representation
 - Calculated and measured properties of chemical compounds are represented in a uniform way
 - Linked to the resource used for data generation
 - Annotated via ontology entries
 - Model representations link to algorithms and data used

All OpenTox components are defined by
OWL ontology

<http://opentox.org/api/1.1/opentox.owl>

All resources are subclasses of
`ot:OpenToxResource`



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.

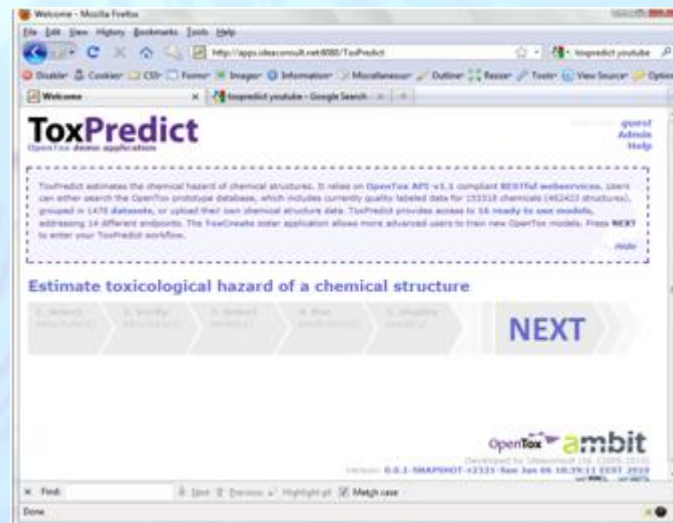
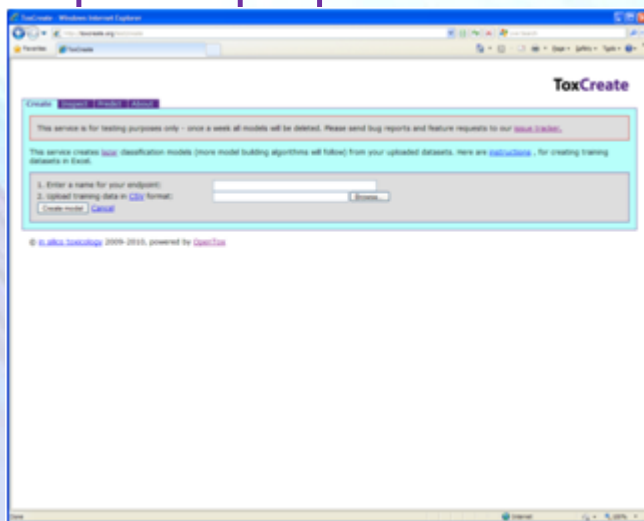
(Subset of) services could be hosted by the same provider, or by multiple providers on separate locations.



Partner No. /Service type	Compound	Dataset	Feature	Algorithm (processing)	Algorithm (model)	Model	Validation	Report	Task	Autherntication and Authorisation service	Ontology service
2	Y	Y		Y	Y	Y			Y		
3	Y	Y	Y	Y	Y	Y			Y		Y
5				Y	Y	Y					
6							Y	Y	Y	Y	
7					Y	Y			Y		
10						Y					

Demo applications

- Two end user oriented demo applications, making use of OpenTox webservices, have been developed, deployed and are available for testing - <http://toxcreate.org> and <http://toxpredict.org> ;
- ToxCreat creates models from user supplied datasets;
- ToxPredict uses existing OpenTox models to estimate chemical compound properties



ToxPredict: a detailed case study

- ToxPredict estimates the chemical hazard of chemical structures. It relies on [OpenTox API-v1.1](#) compliant RESTful webservices. Users can either search the OpenTox prototype database, which includes currently quality labelled data for **~150,000 chemicals**, grouped in datasets, or upload their own chemical structure data. ToxPredict provides access to **14 ready to use models**, addressing **13 different endpoints**
- ToxPredict uses the following OpenTox webservices: [Compound](#), [Feature](#), [Dataset](#), [Algorithm](#), [Model](#), [Task](#) and [Ontology](#);
- more details on its interactions with OpenTox webservices which are taking place behind the scenes and without requiring any end-user intervention;

ToxPredict: Step 1 (Select structure(s))



Find structure by name, registry number, SMILES, InChI, structure, substructure, similarity...



ToxPredict
Web
Application

OT Dataset API *HTTP GET*

OT Dataset
Service

text/uri-list,
application/rdf+xml,
chemical/x-daylight-smiles
chemical/x-mdl-sdfile,...



Here is the list of structures as
URI links, RDF, MOL or SMILES.

[illegible]

Select and/or edit structure(s)

OT Dataset API *HTTP GET*

OT Dataset
Service

ToxPredict
Web
Application

text/uri-list,

Here is the list of structures as URI links, RDF, MOL, SMILES or images.



ToxPredict: Step 3 (Select model(s))



What prediction models are available? Is there a model for endpoint X?

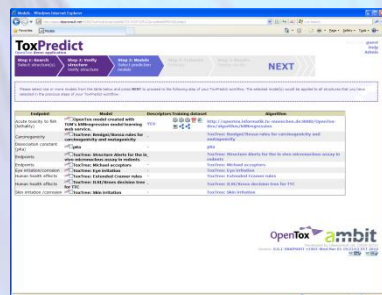


ToxPredict
Web
Application

HTTP GET SPARQL query

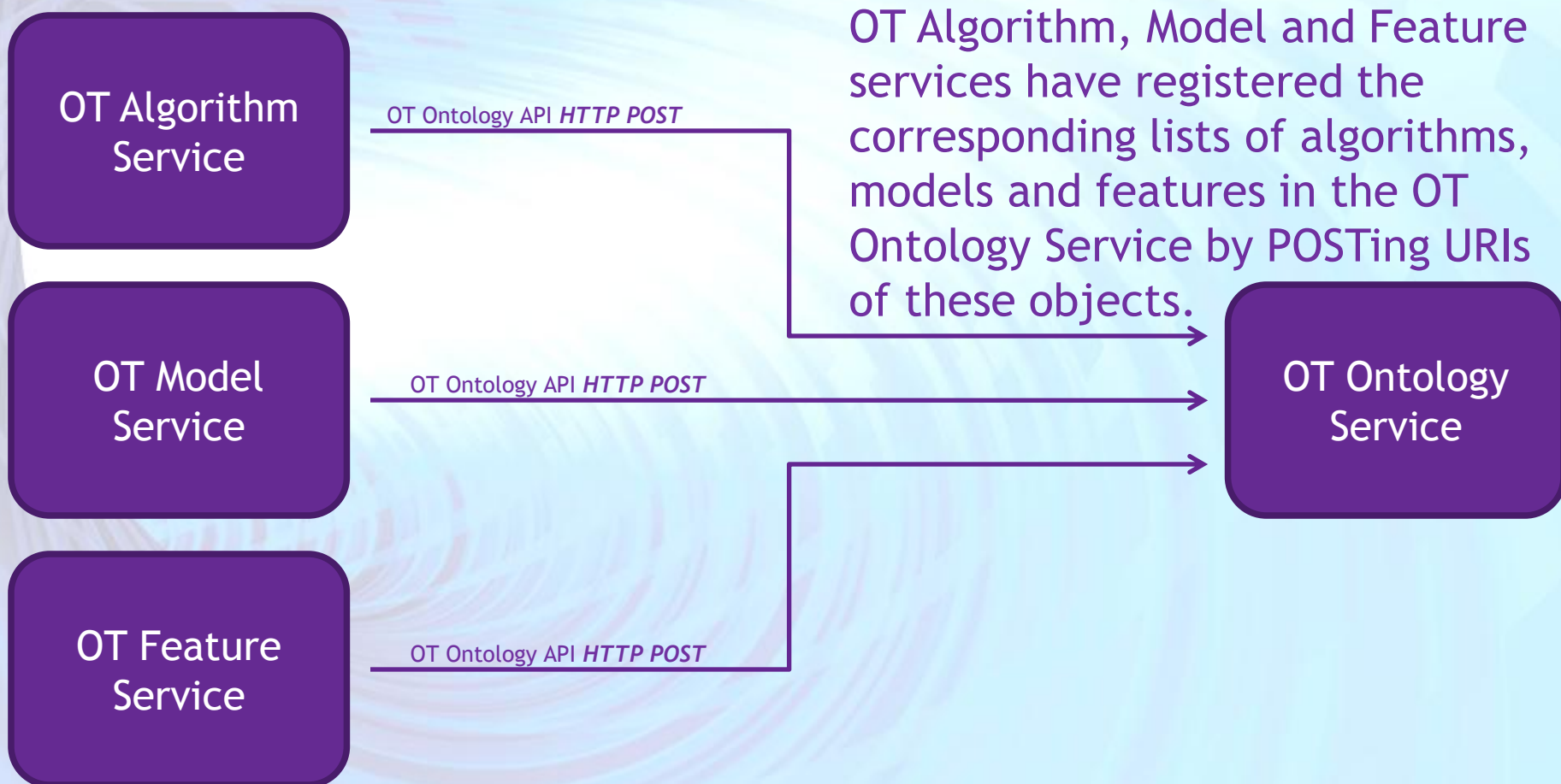
OT Ontology
Service

application/sparql-results+xml

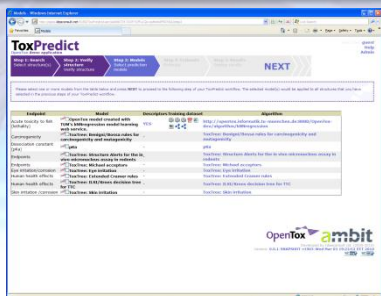


Here is the list of model URIs and related endpoints and algorithms in SPARQL format.

ToxPredict: Step 3 (behind the scenes)



ToxPredict: Step 4 (Estimate)



Run the selected models.

OT Model API *HTTP POST* with
parameter dataset URI from
steps 1-2

OT Model
Service

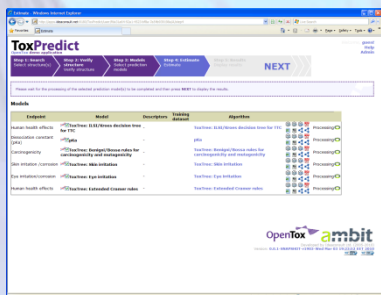
ToxPredict
Web
Application

HTTP code 202 “Accepted”
Model task URI in
HTTP:location header

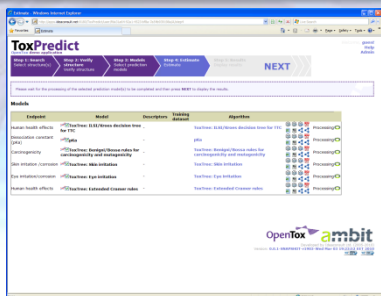
The calculation will take a while;
here is a task URI, which can be
queried for processing status.

Create a
new task.

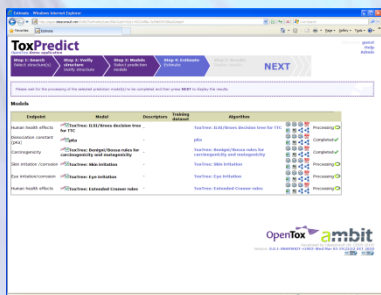
OT Task
Service



ToxPredict: Step 4 (Estimate)



ToxPredict
Web
Application



Is the task completed?

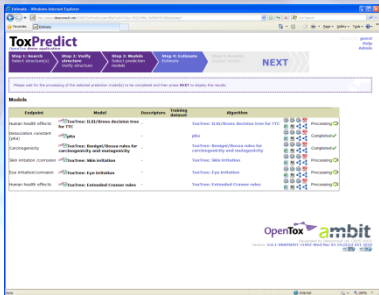
OT Task API *HTTP GET* on task URI

OT Task
Service

HTTP code 202 "Accepted for processing"
Returns Task URI

Not yet, but processing has finished and the results have been posted to the OT Dataset service; here is a task URI for the dataset import.

ToxPredict: Step 4 (Estimate)

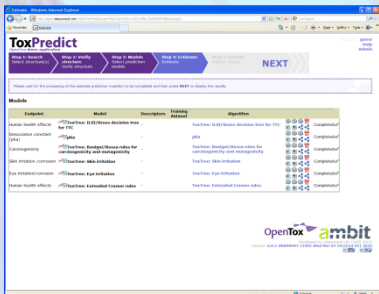


ToxPredict
Web
Application

Is the task completed?

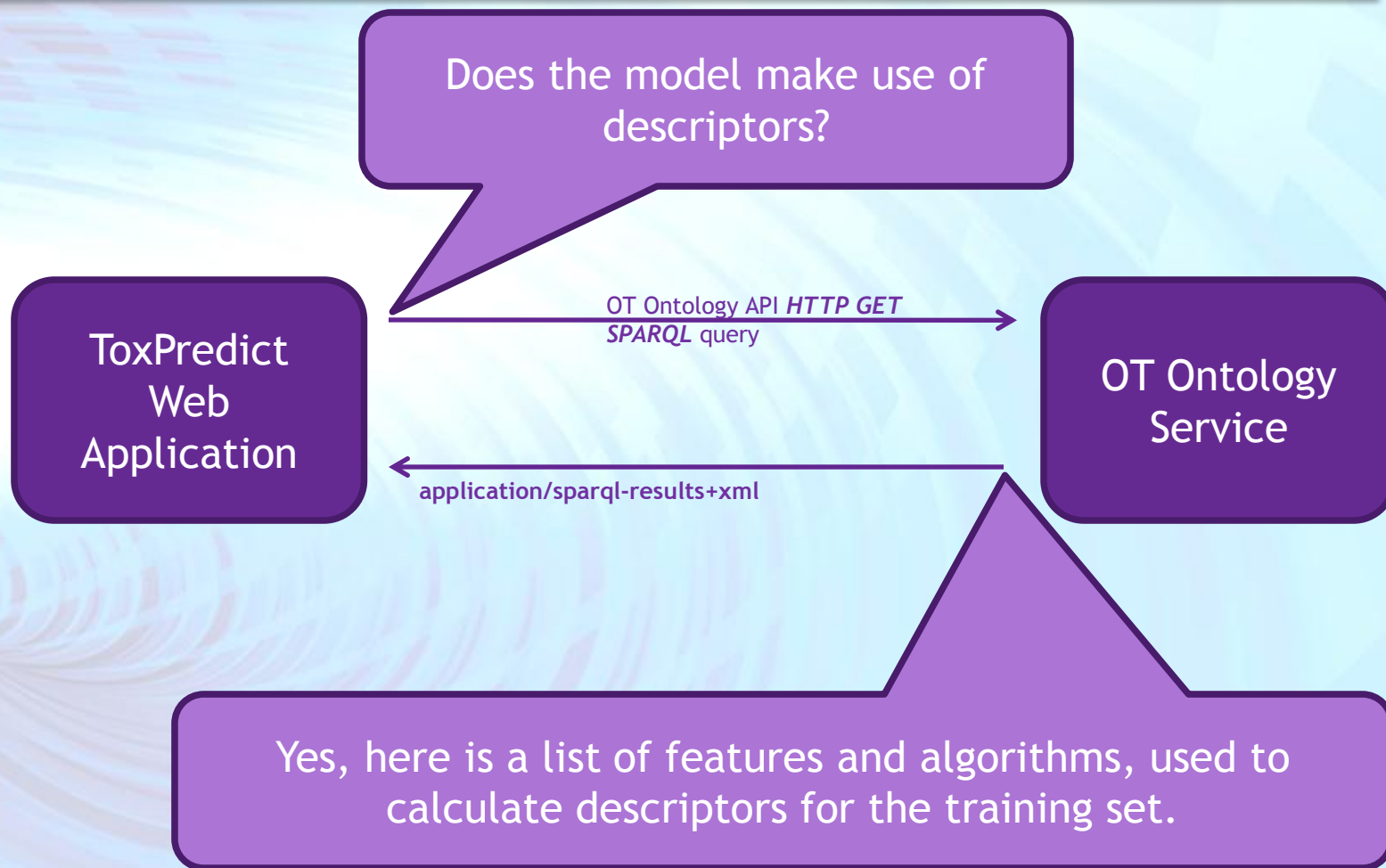
OT Task API *HTTP GET* on task URI

OT Task Service

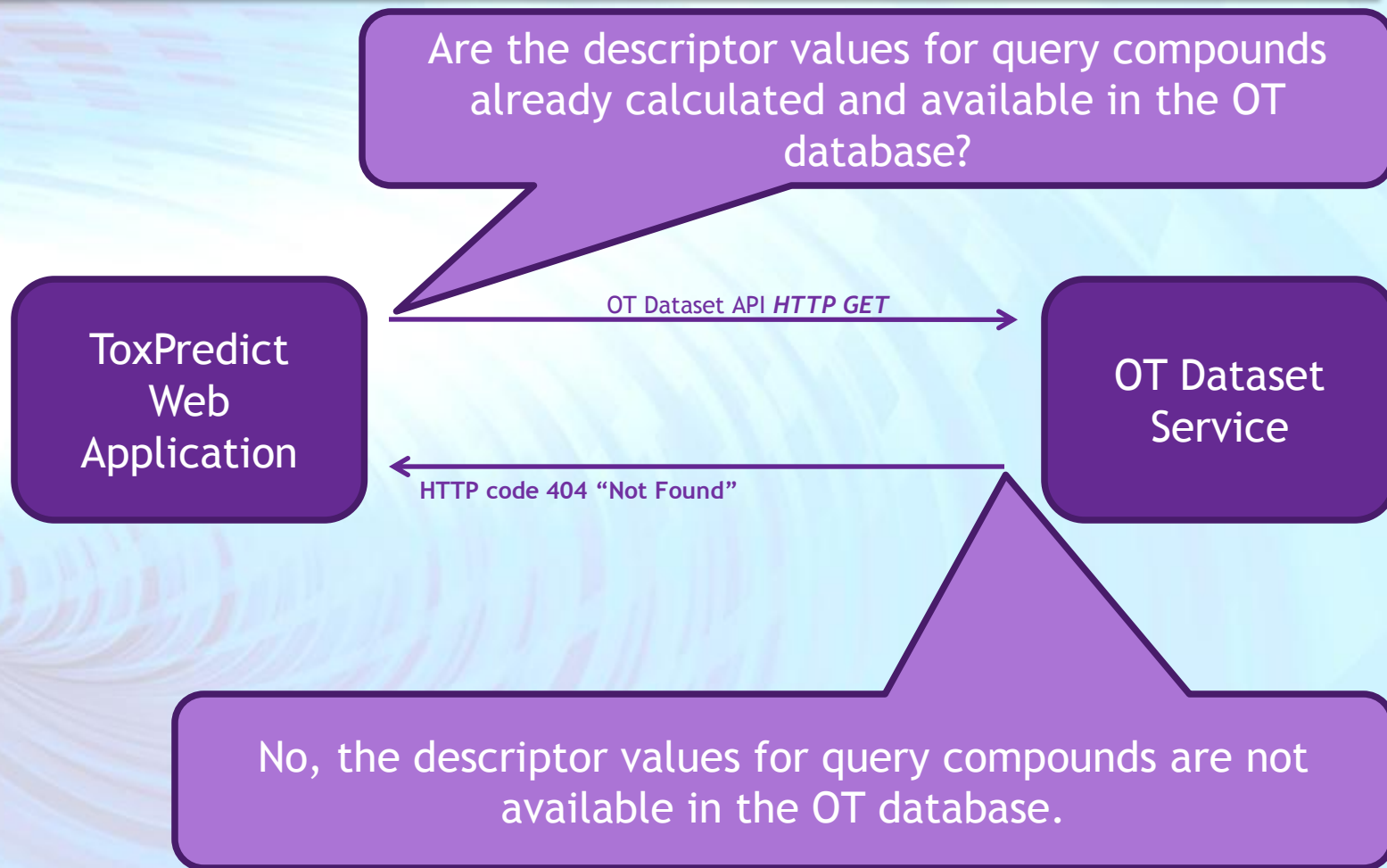


Yes, here is the dataset URI of the results.

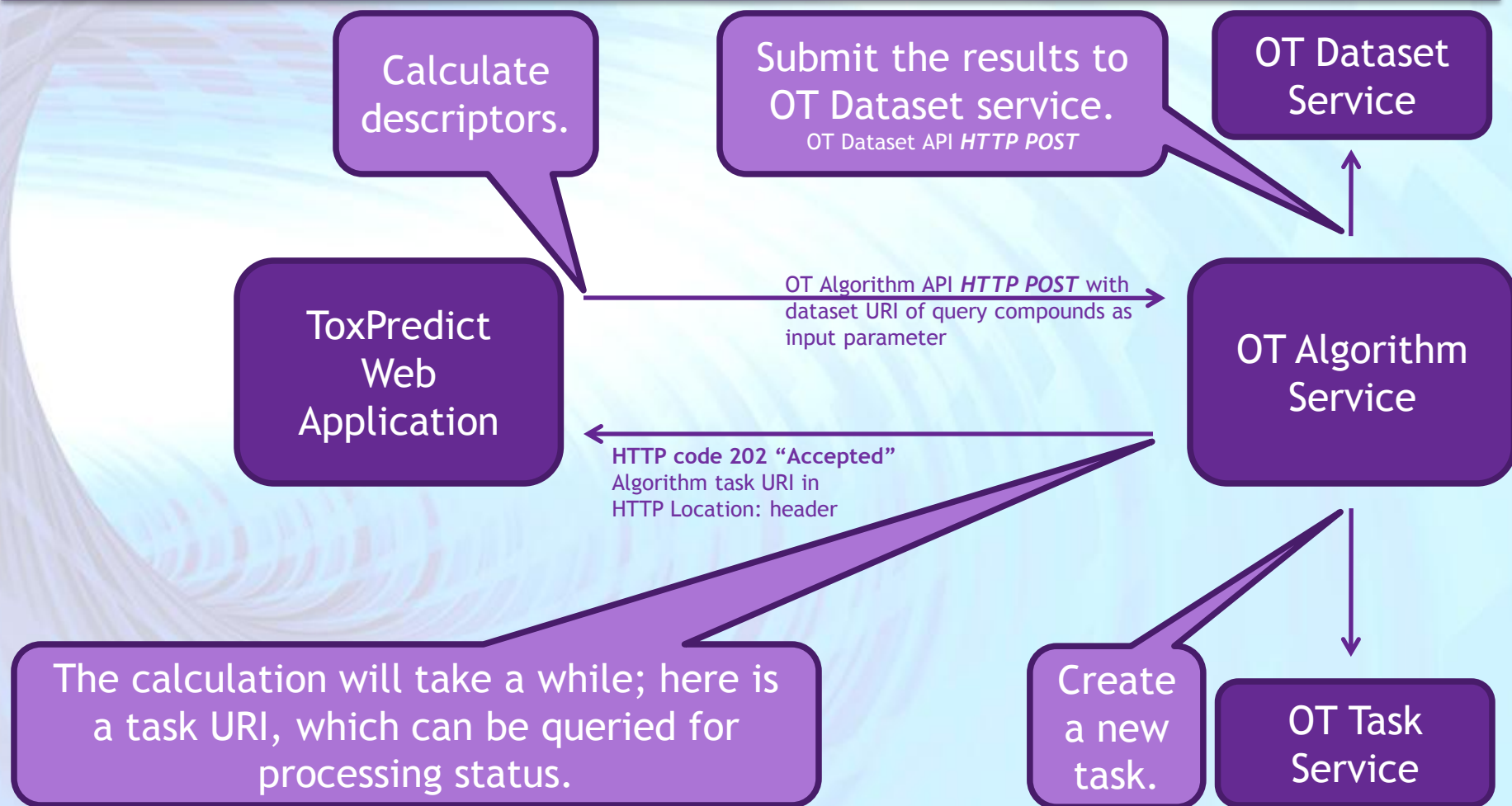
ToxPredict: Step 4 (behind the scenes)



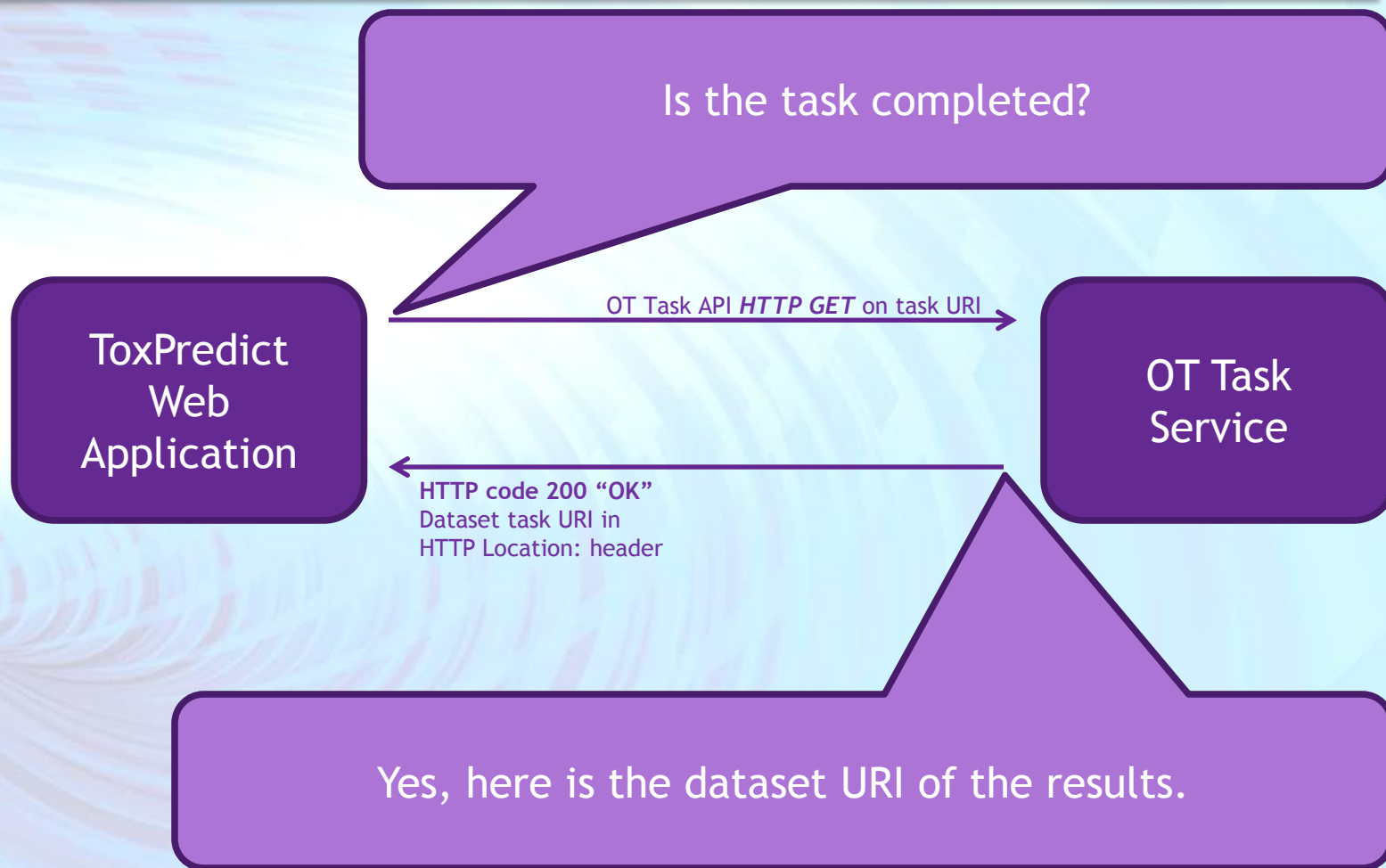
ToxPredict: Step 4 (behind the scenes)



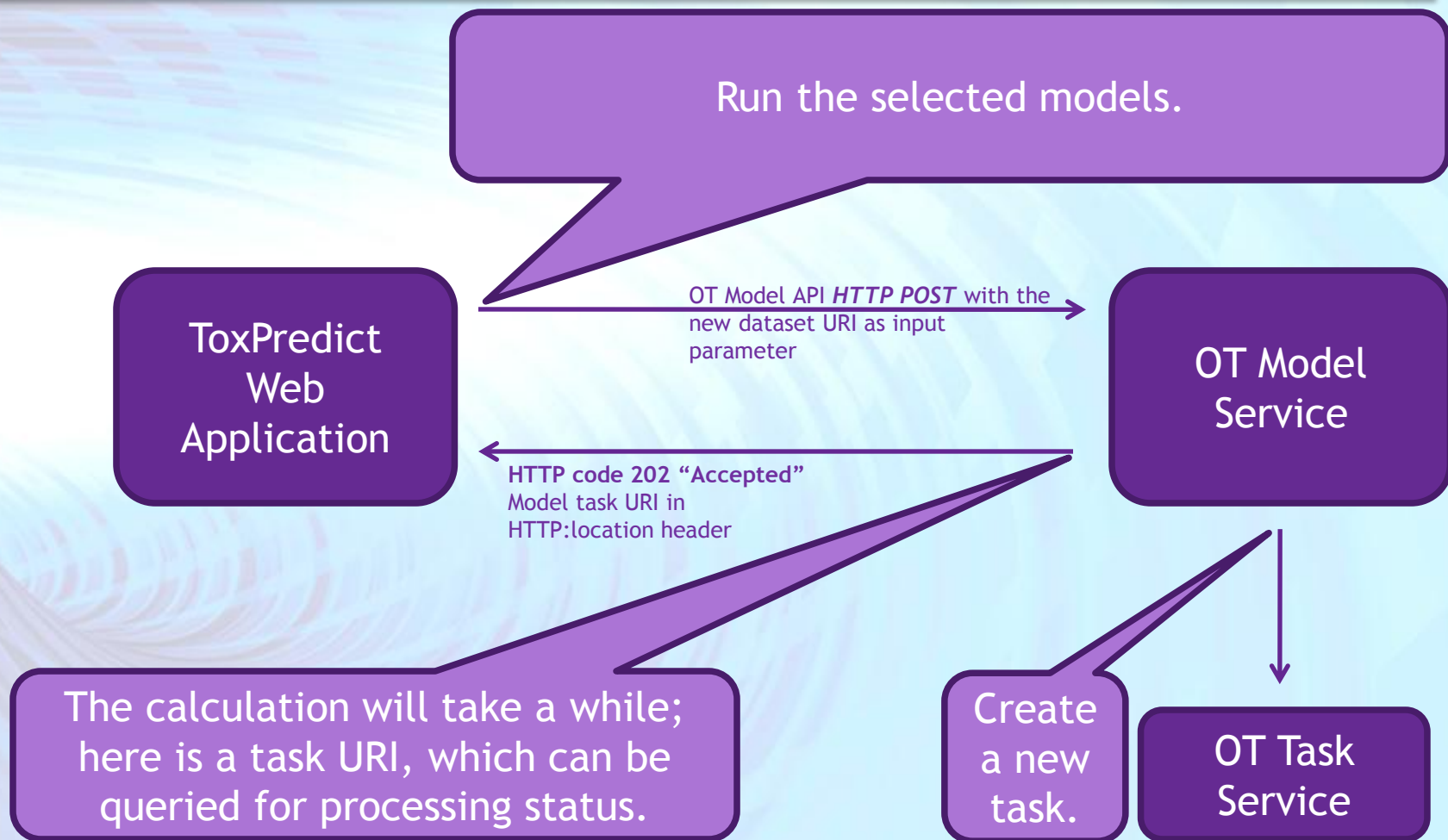
ToxPredict: Step 4 (behind the scenes)



ToxPredict: Step 4 (behind the scenes)

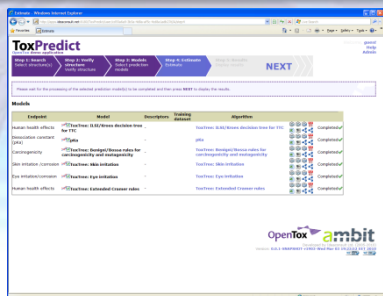


ToxPredict: Step 4 (behind the scenes)



... the remaining part of the Estimation step has been described

ToxPredict: Step 5 (Display results)



Retrieve calculation results from the final dataset URI, obtained in Step 4 (Estimation).

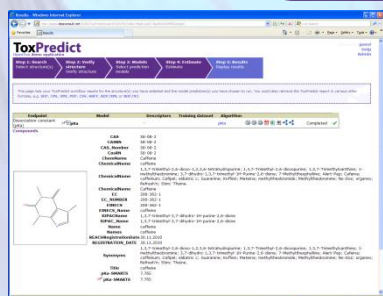


ToxPredict
Web
Application

OT Dataset API *HTTP GET*

OT Dataset
Service

application/rdf+xml



Here is the dataset content in RDF format, according to OpenTox.owl and containing estimation results, as well as compound identifiers and experimental data.

Development and Use of Predictive Toxicology Applications

An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Ontology and Data Schema

Olga Tcheremenskaia
Romualdo Benigni



(Istituto Superiore di Sanità, Rome, Italy)

Ontology: presentation plan

Introduction, Tools, Principles

- What is an ontology?
- OWL and OBO languages, Protégé Editor
- Biomedical ontologies
- OBO Foundry Principles

Incorporation of Ontologies into OpenTox

- Why we need an ontology for the OpenTox?
- Principles of the OpenTox toxicological endpoint ontology development

ToxML data standard

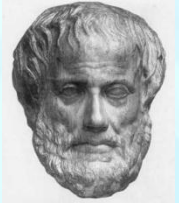
- ToxML data standard
- Toxicological Databases mapping onto ToxML scheme
- ToxML in ontology development

Ontology Development

- Collaborative Protégé Server
- Toxicological Endpoints Ontology structure
- Organs ontology
- OpenToxipedia: community knowledge resource

What Is An Ontology?

Ontologies in Information Science

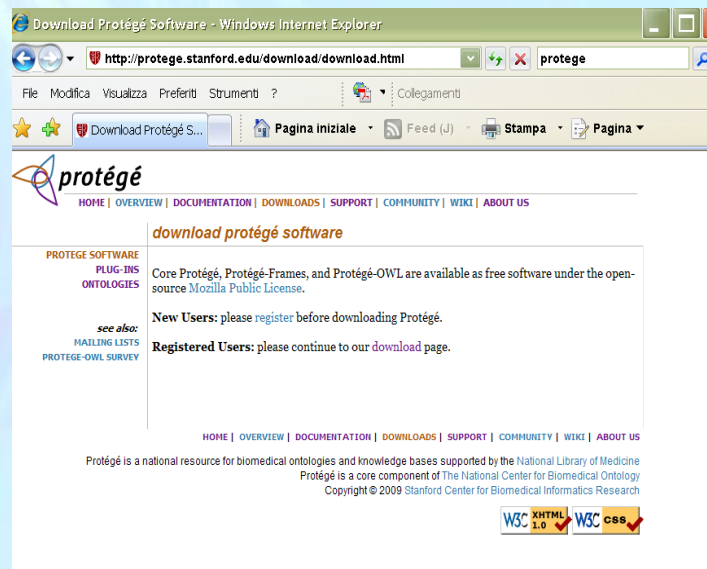


- **Ontology** (Socrates & Aristotle 400-360 BC): The study of being
- Explicit description of the conceptualization of a domain
- **Ontology in Information Science** is a base for **Semantic Web** (a group of methods and technologies to allow machines to understand the meaning - or "semantics" - of information)
- **Ontology Components**
 - **concepts** - set of entities within a domain
 - **relations** - interaction between concepts or concept's properties
 - **instances** - concrete examples of concepts of the domain
 - **axioms** - explicit rules to constrain the use of concepts
- **Role of Ontologies**
 - To share **common understanding** of the structure of descriptive information
 - among people
 - among software agents
 - between people and software
 - To enable **reuse** of domain knowledge
 - To introduce **standards** to allow interoperability

OWL : The Web Ontology Language

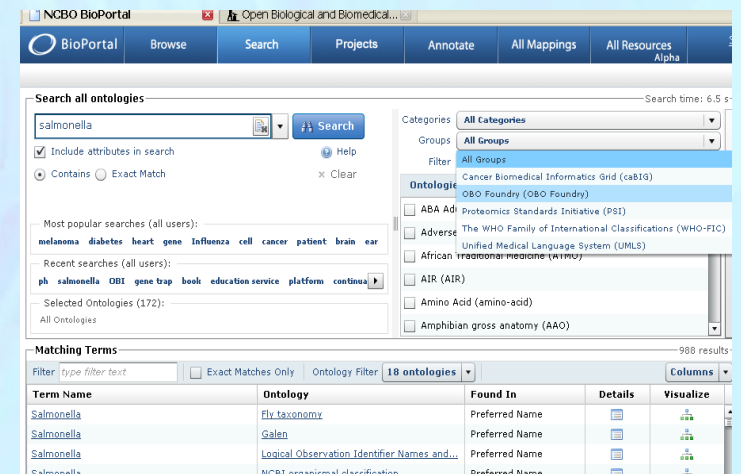
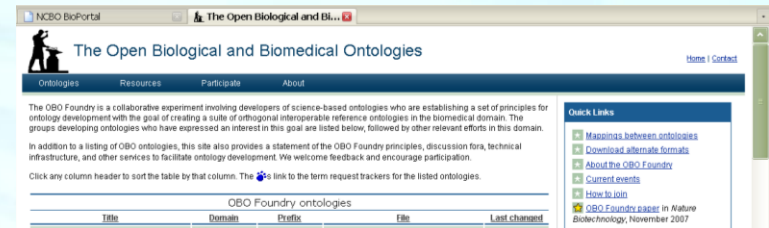
Protégé Ontology Editor

- **OWL:** the Web Ontology Language is the World Wide Web Consortium (W3C) standard
- **OWL components:**
 - a formal description of concepts, terms, and relationships within a given knowledge domain
 - formal computational definitions
 - tools for reasoning
- **OWL tutorial** published by The University of Manchester “A Practical Introduction to Ontologies & OWL” <http://www.co-ode.org/resources/tutorials/intro/>
- **Protégé** is a free, open source ontology editor and knowledgebase framework (<http://protege.stanford.edu>)
 - provides tools for visualizing ontologies as well as for constructing them
 - facilitates OWL ontology development and maintenance
 - allows an automated reasoning: RACER is a reasoner frequently used with Protégé (www.sts.tu-harburg.de/~r.f.moeller/racer/)



Public Biomedical Ontologies Resources

- **The Open Biomedical Ontologies (OBO):** ontology development is regulated within the OBO Foundry, which defines a set of shared principles governing ontology development
<http://obofoundry.org>
- **OBO Foundry principles and criteria:**
 - The ontology is **open** and available to be used by all.
 - The ontology is in a **common formal language**.
 - The developers of the ontology agree in advance to **collaborate** with developers of other OBO Foundry ontology **where domains overlap**.
 - **Orthogonality:** for any particular domain, there is community convergence on a single controlled vocabulary.
- **Web-based ontology portal: BioPortal** (www.bioontology.org/tools/portal/bioportal.html) allow users to browse, search, and visualize ontologies.
- **Ontology Lookup Service:** accessing ontologies www.ebi.ac.uk/ontology-lookup/

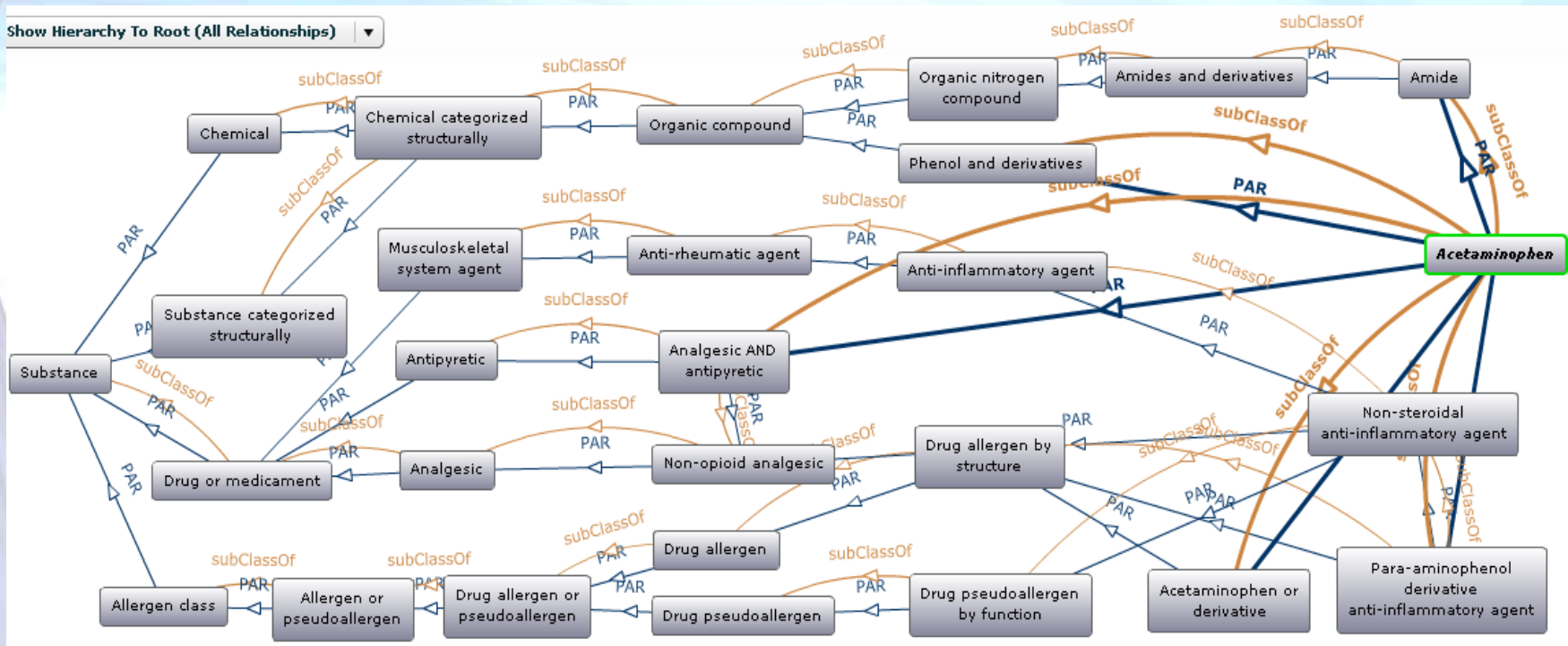


Biomedical Application of Ontologies and Examples

- **Possible biomedical applications of ontologies:**
 - Search and query of heterogeneous biomedical data
 - Data exchange among applications
 - Information integration
 - Representation of encyclopedic knowledge
 - Computer reasoning with data
- **A growing community interest in using and producing biomedical ontologies.**
- **Examples of most important ontology projects:**
 - Gene Ontology (GO)
 - Chemical entities of biological interest (ChEBI)
 - Ontology for biomedical investigations (OBI)
 - NCI Thesaurus
 - Foundational Model of Anatomy (FMA)
 - Mouse Adult Gross Anatomy (MA)
 - Mouse gross anatomy and development (EMAP)
- **Toxicology domain - no ontology available**

Example of ontological representation of terms

SNOMED CT (Systematized Nomenclature of Medicine - Clinical Terms):
Paracetamol (Acetaminophen) example



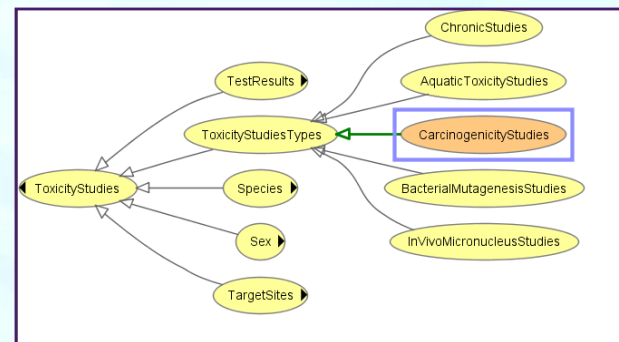
OpenTox Toxicological Endpoint Ontology

•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



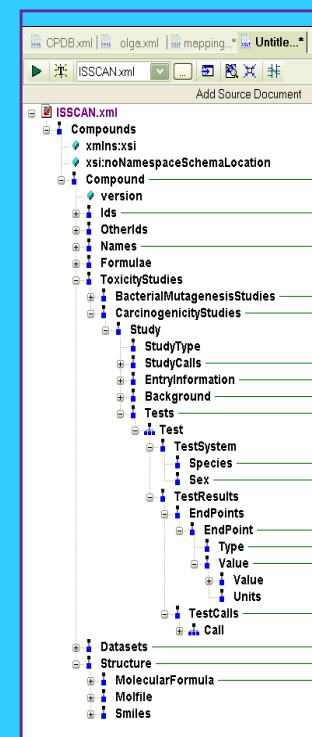
Other freely available resources: DSSTOX, GoReni (ITEM), etc

Protégé, free open source OWL (Web Ontology Language) editor

OpenTox
Toxicological
Data Ontology

Re-use of pieces or terms defined in neighboring ontologies (OWL and OBO)

ToxML scheme



Toxicological data: needs for standards

ToxML scheme

- Need for data standards for automatic data integration
 - Example: **Carcinogenic Activity**

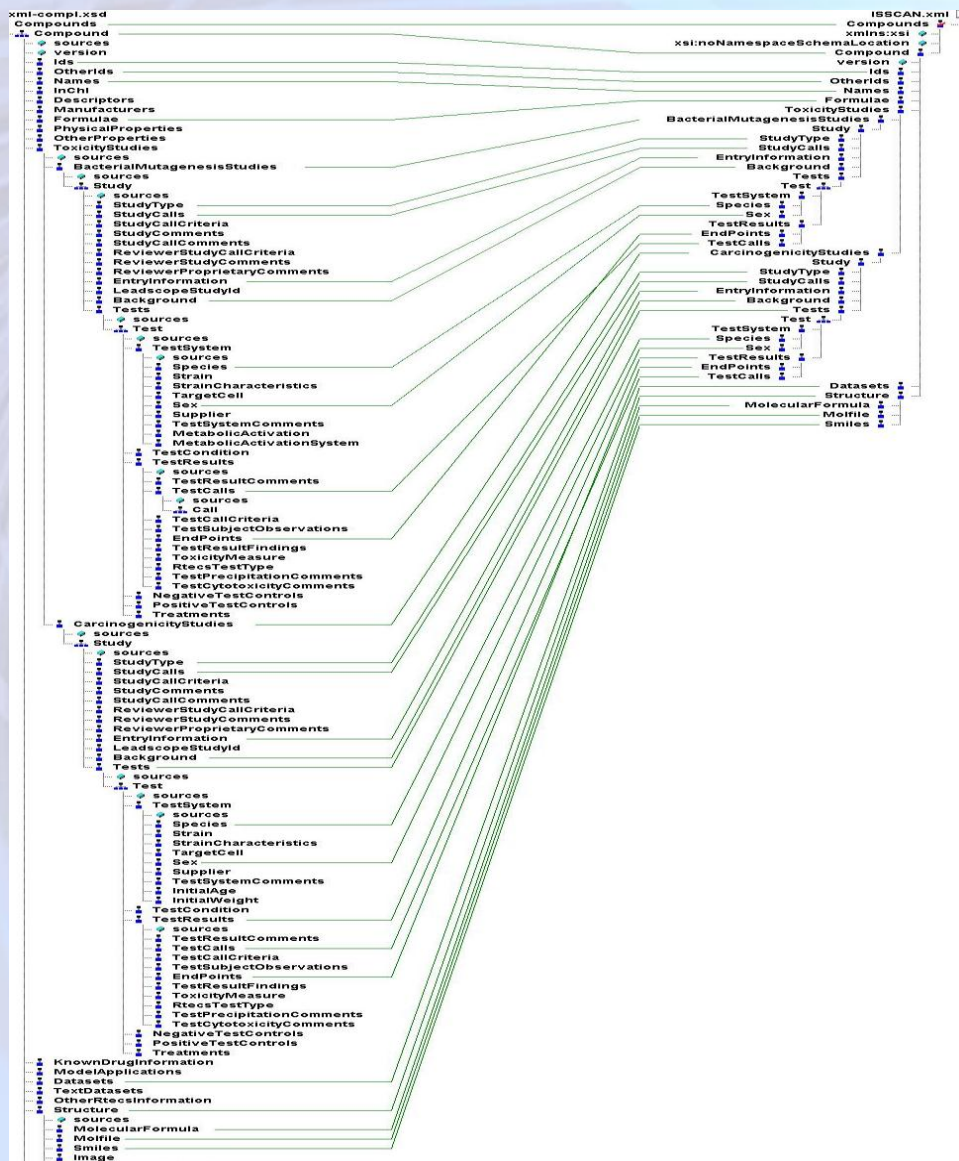
CPDBAS: Carcinogenic Potency Database
http://www.epa.gov/ncct/dsstox/sdf_cpdbas.html#SDFFields

ISSCAN: Chemical Carcinogens Database
<http://www.iss.it/ampp/dati/cont.php?id=233&lang=1&tipo=7>



- **ToxML** (<http://www.leadscope.com/toxml.php>)
 - is a public initiative led by scientists at LeadScope, Inc
 - controlled vocabularies and XML scheme for storing chemical toxicity data
 - the latest version of ToxML public schema (April 7th , 2009)
- It is supported by OpenTox for interoperable data communications between services

Mapping of the ISSCAN entry - ToxML xsd scheme

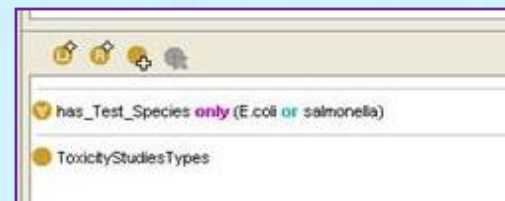


CAS	67-66-3
Substance ID	4
ChemName	Chloroform Formyl trichloride; methane trichloride; methenyl trichloride; Methyl trichloride; R 20; r 20 (refrigerant); Refrigerant R20; trichloroform; Trichloromethane
Synonyms	
SAL	negative
Canc Reference	positive CPDB
TD50_Rat	262
TD50_Mouse	90.3
Rat_Male_Canc	positive
Rat_Male_NTP	ND
Rat_Female_Canc	positive
Rat_Female_NTP	ND
Mouse_Male_Canc	positive
Mouse_Male_NTP	ND
Mouse_Female_Canc	positive
Mouse_Female_NTP	ND
MolWeight	119.38
Formula	CHCl3
SDF file	Connection table
SMILES	C1C(Cl)Cl

ToxML: mapping of toxicological data

Needs for extensions

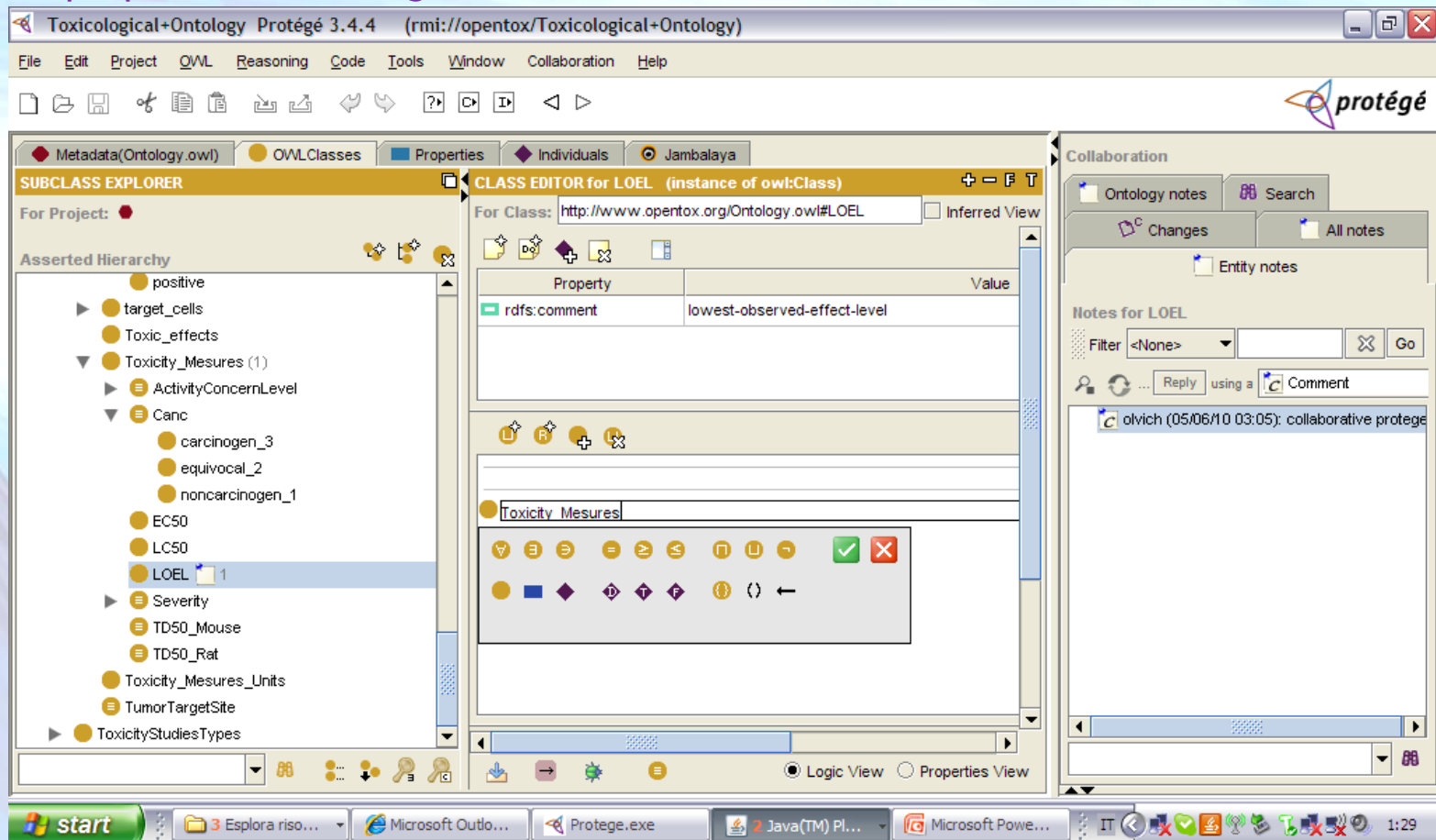
- The latest version of public ToxML scheme (Last updated April 7th, 2009) has been studied in terms of their suitability to map the content of databases candidate for the OpenTox database using Stylus Studio 2008 XML Enterprise software
 - 1. ISSCAN database
 - 2. In vivo micronucleus database (in development at the ISS);
 - 3. Bacterial mutagenesis database (to be developed at the ISS)
 - 4. RepDose ITEM database (ITEM)
- **OpenTox Toxicological Ontology: ToxML integration**
 - OWL Ontology: Classes and Hierarchies on the base of ToxML scheme
 - Extension needs:
 - Free text field: TargetSite - Organs Ontology and Effects Ontology to avoid a huge variability of terms, as each laboratory/researcher is able to enter his individual description of an observed effect
 - Introduction of relationship between different classes, restrictions: e.g., introducing the property “has_Test_Species” limits the test species suitable for certain toxicity study.



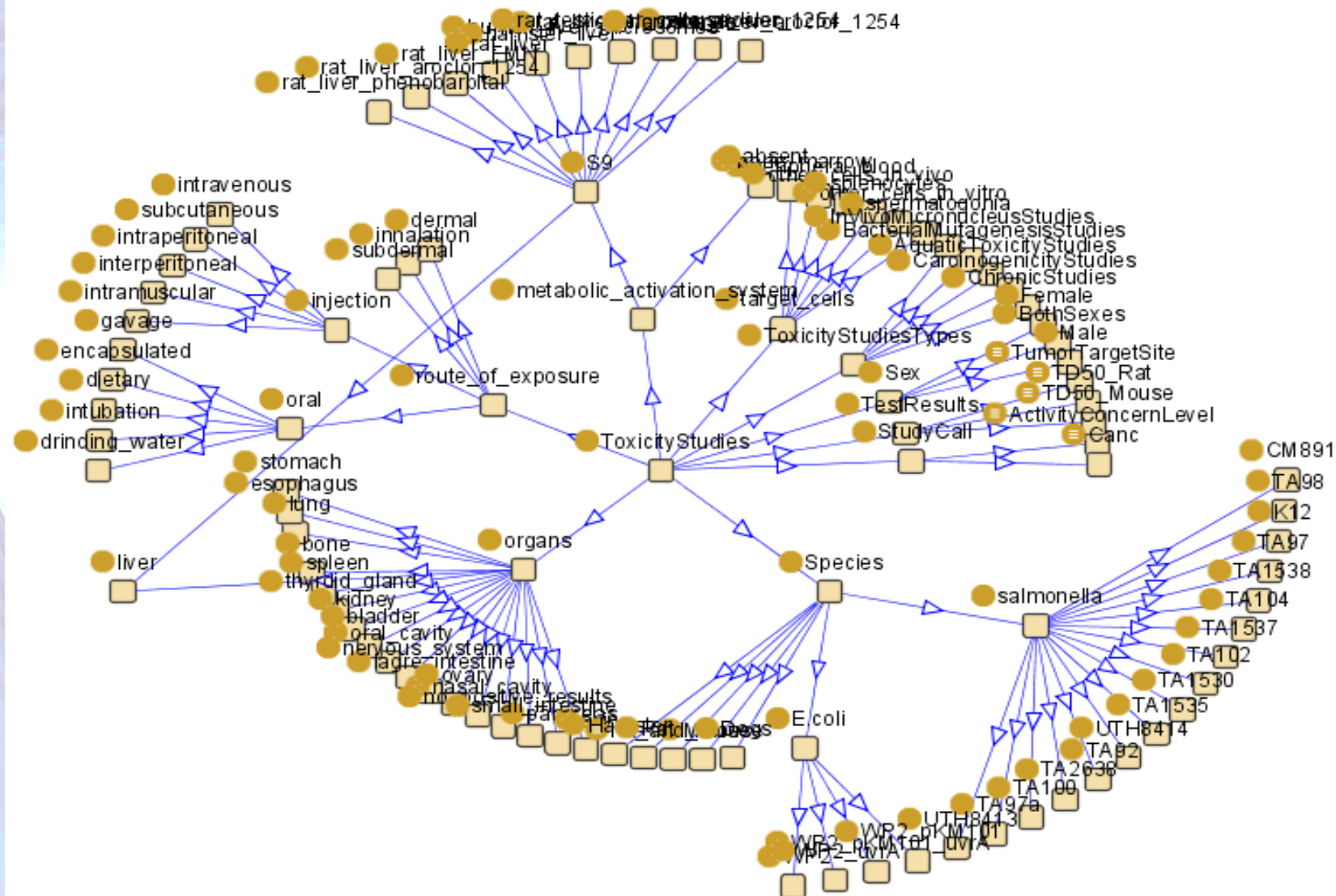
OpenTox: Collaborative Ontology Development

- **Collaborative Protégé: why?**

- annotating ontology components
- tracking changes, history of concept
- discussion, live messages
- proposals and voting
- searching and filtering
- defining users, groups, policies
- available in multi-user mode



Toxicological Ontology: graphical representation



Ontology For Target Organs

- Contribution of Fraunhofer Institute for Toxicology and Experimental Medicine (ITEM)
- Ontology Development on the base of INHAND (International Harmonization of Nomenclature and Diagnostic Criteria for Lesions in Rats and Mice)

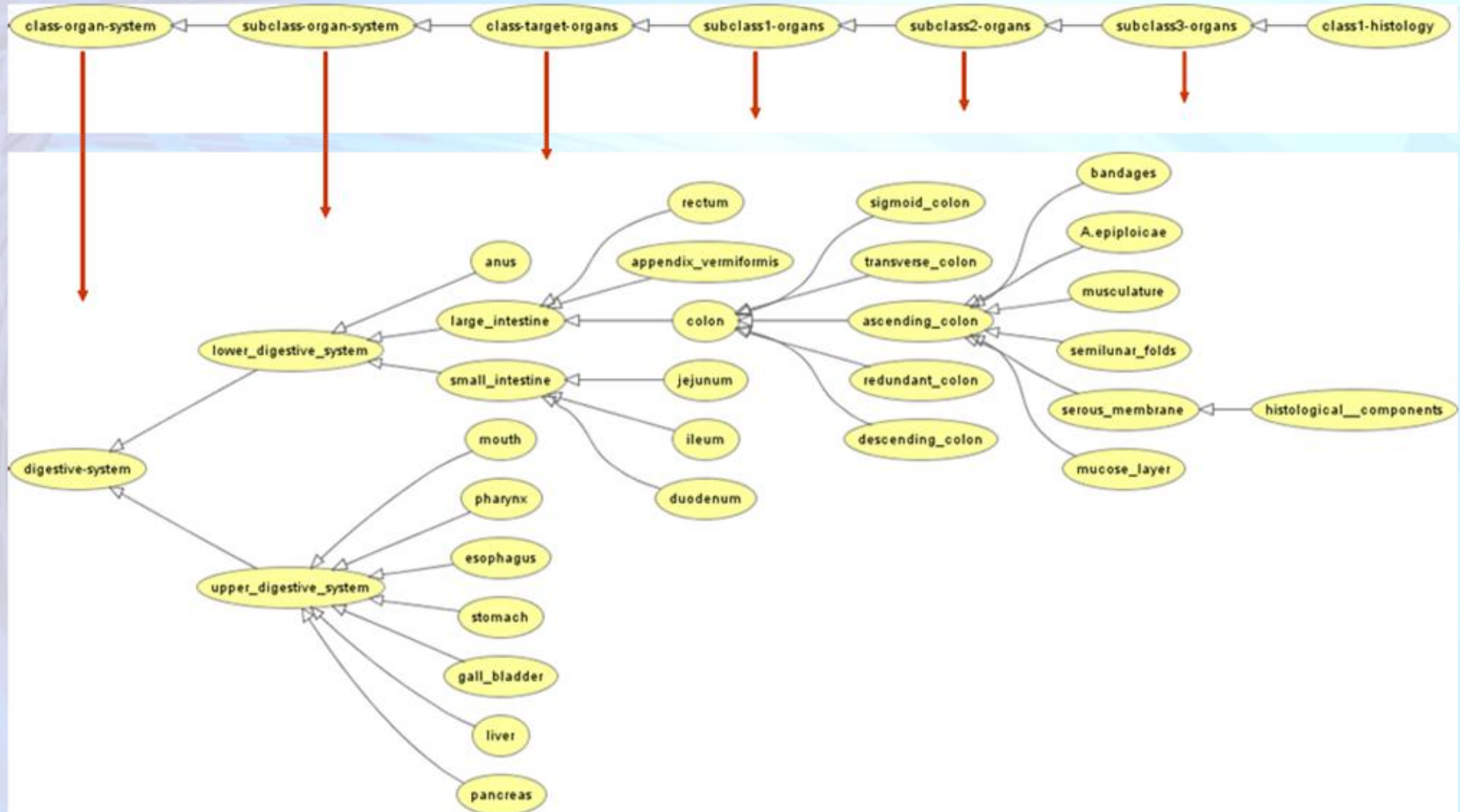
Organs system (class) - Subclass organs system

Target organs (class) - Targets organs (subclass 1 to N)

Histopathology (class) - Histopathology (subclasses if needed)

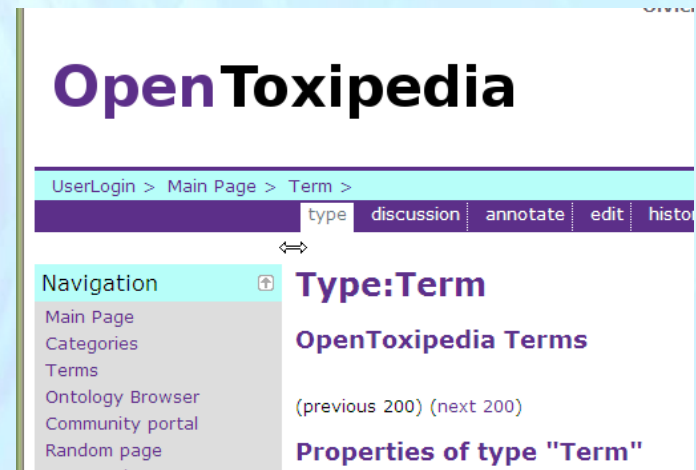
- **Key Features of the ITEM Organs Ontology:**
 - more synonyms regarding organ systems, target organs and their subclasses
 - organs are more detailed, up to histological components
 - organs linked to different organ systems (e.g. ovaries)
 - reviewed by pathologists, who have been involved in the INHAND process
 - 12 organ systems fully described
- **Perspectives:**
 - linking of the organ systems and their components with pathologic effects

Ontology For Target Organs: Digestive system example



OpenToxipedia: community-based, predictive toxicology knowledge resource

- The OpenToxipedia content has moved to the new Semantic Media Wiki (SMW) platform. You will find a link to it at the common page <http://www.opentoxipedia.org>
 - Community based collaborative database
 - Automatically-generated lists.
 - Improved data structure, semantic terms annotation
 - External reuse: RDF export and import in Protégé Ontology Editor
- Creating, adding, editing and keeping terms used in toxicology terminology
- 862 toxicological terms with description and literature references classified into 26 categories
- Transparency and scientific basis of information
- Curation by the OpenTox toxicologists



Plan for Future Work

- Continue collaborative ontology development using the Collaborative Protégé server, covering more toxicological endpoints
- Review carefully all existing OBO and Bioportal ontologies, find the overlapping ontology, valuate the possibility of import
- External collaboration, e.g., with Ontology for Biomedical Investigations (OBI) group
- Setup Protege sub project for *in vitro* assays ontology (ToxCast data)
- Upload OpenTox Toxicological Endpoints Ontology to the BioPortal Web Site
- OpenToxipedia: community based development

OpenTox Ontology Working Group

- **External collaboration**
 - Barry Hardy
- **Algorithms and features ontology**
 - Nina Jeliaskova, Vedrin Jeliaskov, Ivelina Nikolova
 - Christoph Helma
 - Tobias Girschick
 - Andreas Karwath
 - Georgia Melagraki
 - Sunil Chawla
 - David Gallagher
- **Collaborative Protégé administrator**
 - Micha Rautenberg
- **Toxicological Endpoint Ontology**
 - Aldo Benigni, Sylvia Escher, Helvi Grimm, Alexey Lagunin, Olga Tcheremenskaia
- **OpenToxipedia**
 - Alexey Lagunin, Sergey Novikov, Natalya Skvortsova

Development and Use of Predictive Toxicology Applications

An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Data management and integration

presented by Nina Jeliaskova
(Ideaconsult Ltd., Bulgaria)

Outline

- Ontology Server
- Using the Dataset API
- Dataset Integration
- Data quality and accuracy

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/{modelid}
Algorithm	OpenTox algorithm services	http://host:port/algorithm/{algorithmid}
Validation, Report	A validation corresponds to the validation of a model on a test dataset.	http://host:port/validation/{validationid} 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

- Distributed Web services for predictive toxicology
- REST technology
 - Every object has an unique URI
 - URIs are dereferensable
 - Multiple representation of an object is encouraged (e.g. RDF, but also others)
 - Fixed operations - GET, PUT, POST, DELETE
- Every object has RDF representation
 - Compounds
 - Datasets
 - Compound properties
 - Prediction algorithms
 - Models
 - Validation statistics
 - Reports
- Ontologies: Opentox.owl, Blue Obelisk algorithm ontology, OpenTox algorithm types ontology, OpenTox endpoints ontology, based on ECHA endpoints classification; specific endpoints ontologies, developed by ISS &

Ontology service

- RDF triple storage
- REST interface for registration of OpenTox objects
 - HTTP POST
- SPARQL
- query

Search OpenTox RDF - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://apps.ideaconsult.net:8080/ontology/query/Endpoints

Disable Cookies CSS Forms Images Information Miscellaneous Outline Resize Tools View Source Options

Search OpenTox RDF

Features Algorithms Models Endpoints

Import RDF data into Ontology service

URL

SUBMIT

Ontology service 17969 triples

SPARQL

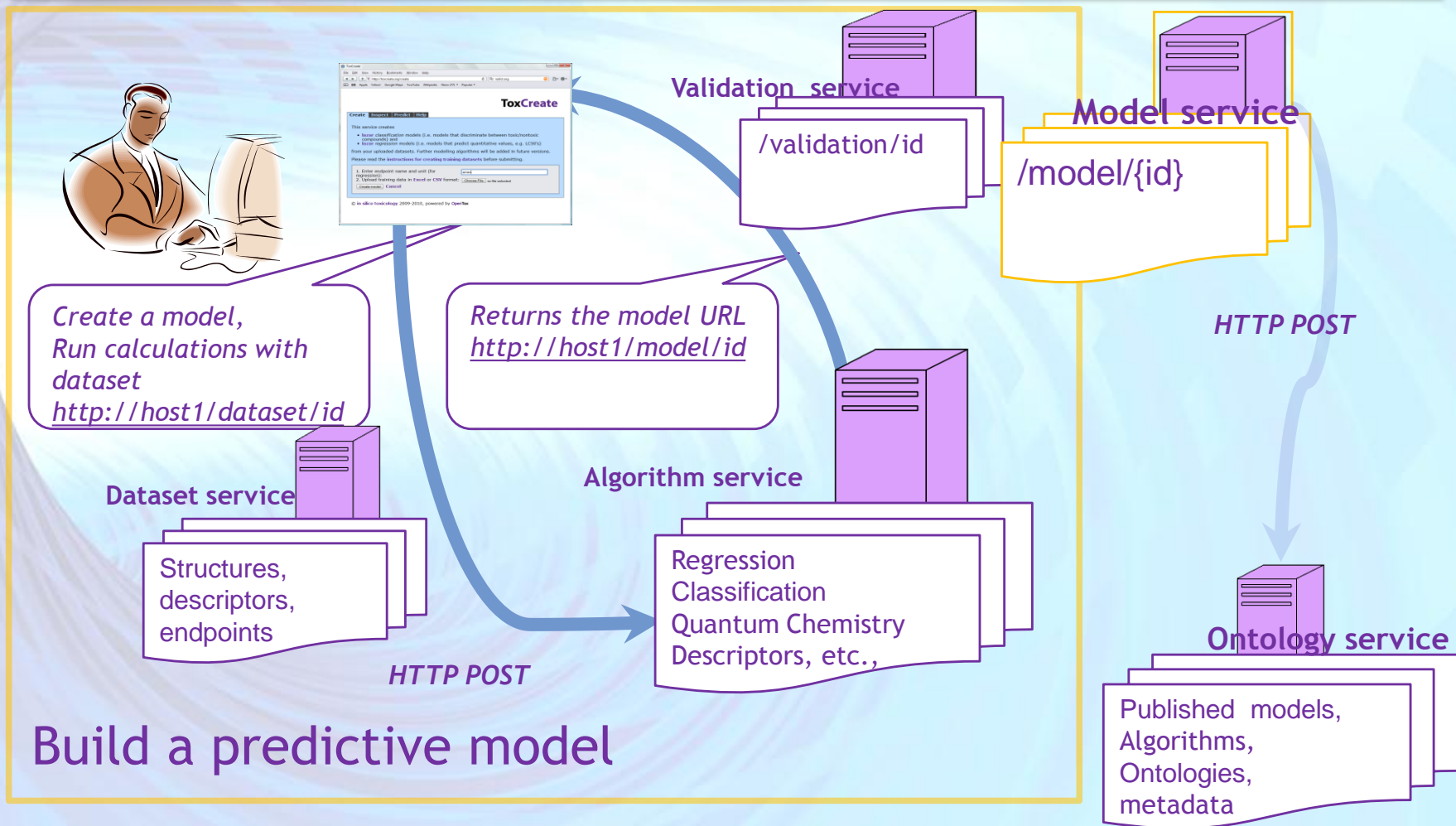
```
PREFIX ot:<http://www.opentox.org/api/1.1#>
PREFIX otee:<http://www.opentox.org/echaEndpoints.owl#>
select ?Endpoints ?title ?id
where {
  ?Endpoints rdfs:subClassOf otee:Endpoints.
  OPTIONAL (?Endpoints dc:title ?title).
  OPTIONAL (?Endpoints dc:identifier ?id).
}
```

Submit Query

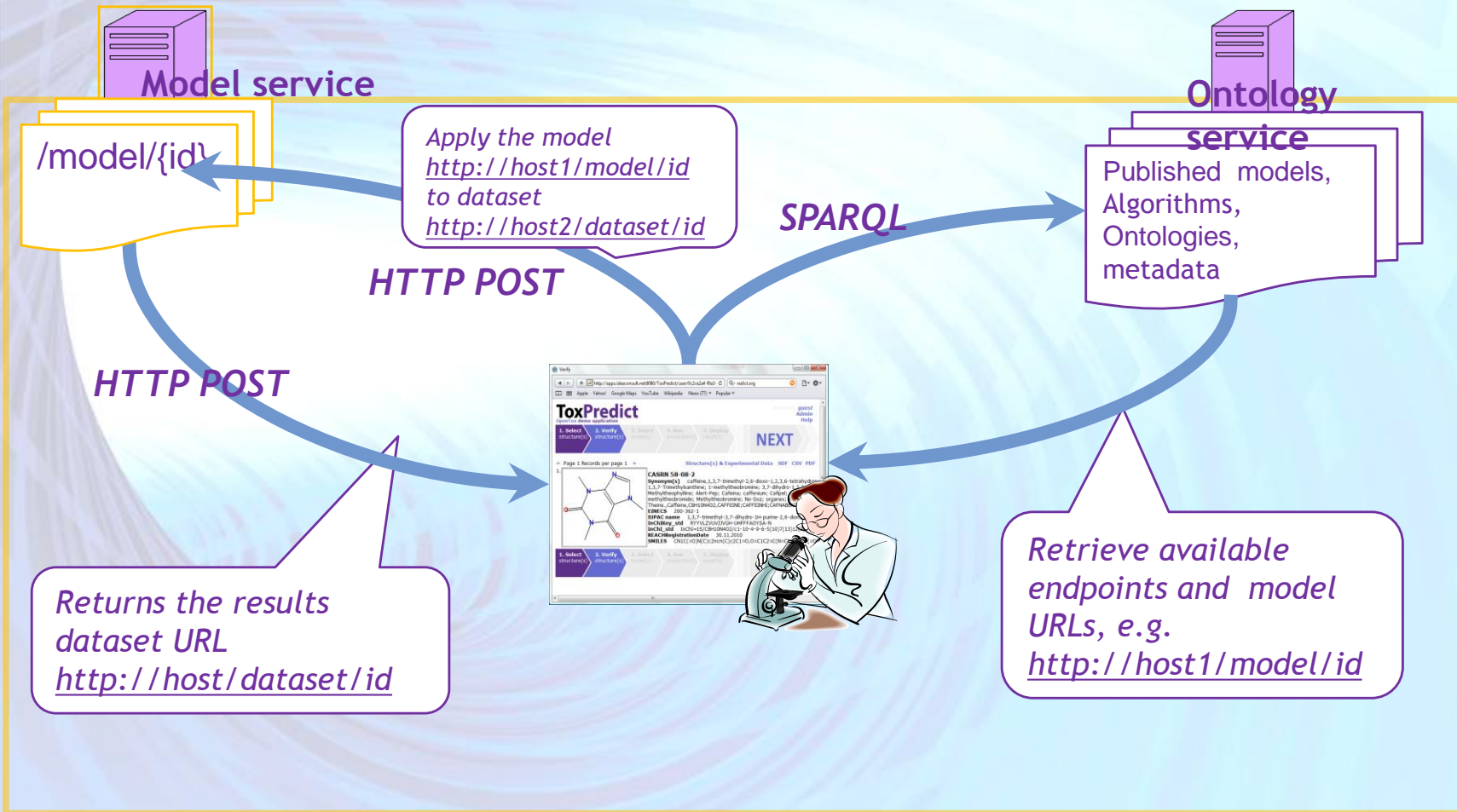
Results [found in 1 ms]

Endpoints	title	id
http://www.opentox.org/echaEndpoints.owl#PhysicoChemicalEffects	Physicochemical effects ""http://www.w3.org/2001/XMLSchema#string	"1""http://www.w3.org/2001/XMLSchema#string
http://www.opentox.org/echaEndpoints.owl#Toxicokinetics	Toxicokinetics ""http://www.w3.org/2001/XMLSchema#string	"5""http://www.w3.org/2001/XMLSchema#string
http://www.opentox.org/echaEndpoints.owl#EcotoxicEffects	Ecotoxic effects""http://www.w3.org/2001/XMLSchema#string	"3""http://www.w3.org/2001/XMLSchema#string
http://www.opentox.org/echaEndpoints.owl#EnvironmentalFateParameters	Environmental fate parameters ""http://www.w3.org/2001/XMLSchema#string	"2""http://www.w3.org/2001/XMLSchema#string
http://www.opentox.org/echaEndpoints.owl#HumanHealthEffects	Human health effects""http://www.w3.org/2001/XMLSchema#string	"4""http://www.w3.org/2001/XMLSchema#string

Build a predictive model

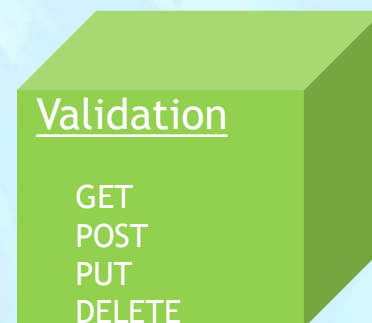
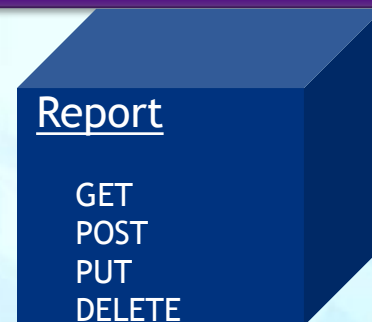
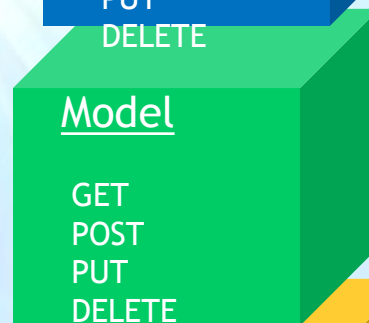
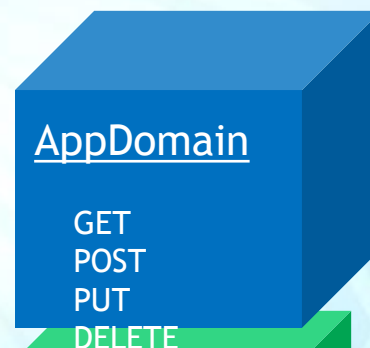
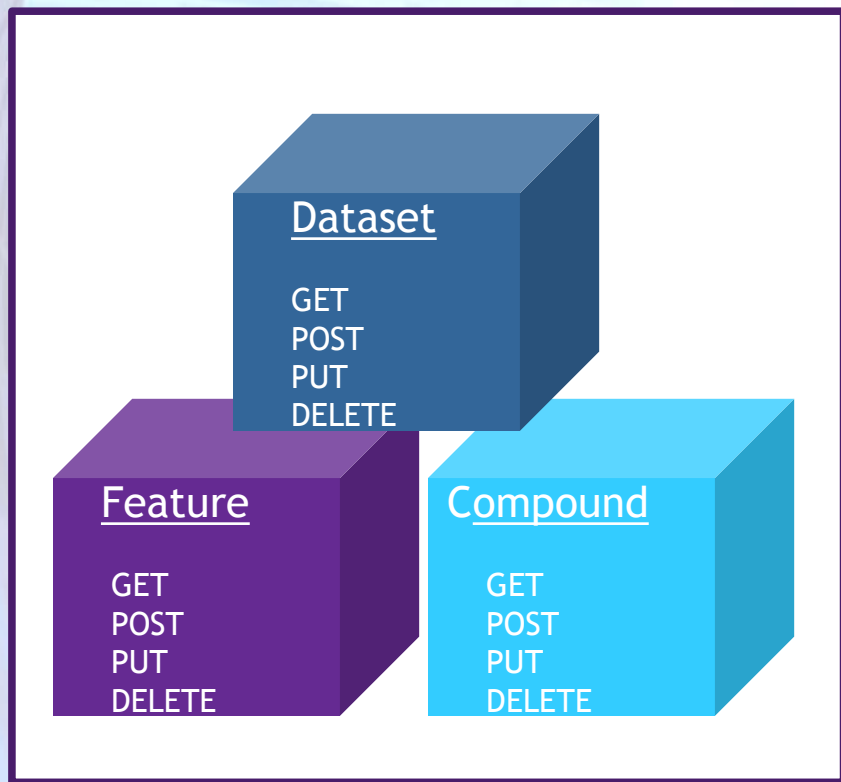


Apply predictive models



OpenTox API (Application Programming Interface)

- The way applications talk to each other
- The way developers talk to applications



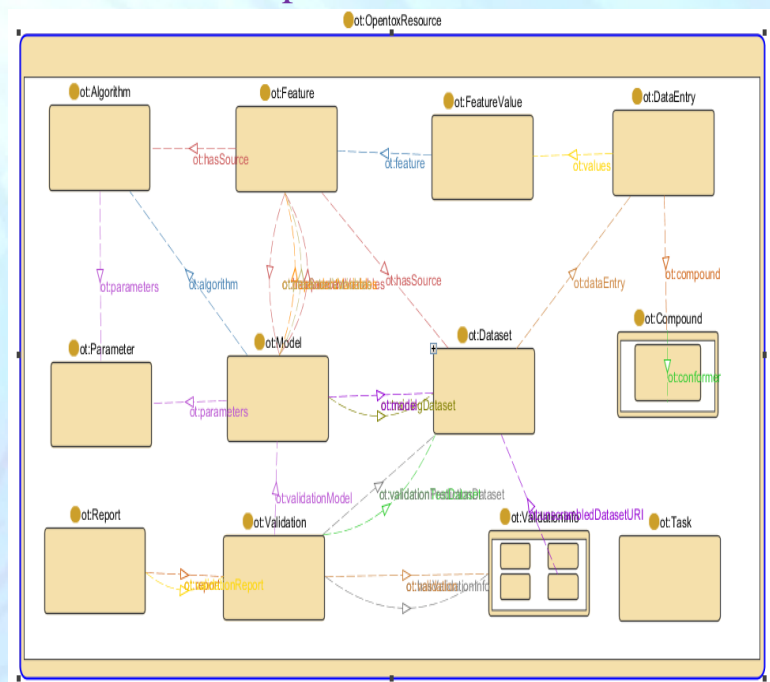
RDF - Resources representation

- The [opentox.owl](http://opentox.org/api/1.1/opentox.owl) ontology
 - A common OWL data model of all OpenTox resources
 - Describes OpenTox resources
 - Describes relationships between them
 - Generates object's RDF representations.
- RDF/XML representation is mandatory for OpenTox resources.
- Uniform approach to data representation
 - Calculated and measured properties of chemical compounds are represented in a uniform way
 - Linked to the resource used for data generation
 - Annotated via ontology entries
 - Model representations link to algorithms and data used

All OpenTox components are defined by
OWL ontology

<http://opentox.org/api/1.1/opentox.owl>

All resources are subclasses of
`ot:OpenToxResource`



Resources: Chemical compound

Compound

Provides different representations for chemical compounds with a unique and defined chemical structure.

`/compound/{id}`

Conformer

`/compound/{id}/conformer/{id}`

Documentation

<http://opentox.org/dev/apis/api-1.1/structure>

Representation

A subclass of `ot:OpenToxResource`.

Supports different Chemical MIME formats

RDF representation only for specifying owl:sameAs links to external resources

Example 1. Retrieve compound as MOL

```
$ curl -H "Accept:chemical/x-mdl-molfile"
http://apps.ideaconsult.net:8080/ambit2/compound/1
CH2O
APtclcactv09040902283D 0 0.00000 0.00000
 4 3 0 0 0 0 0 0 0 0999 V2000
-0.6004 0.0000 0.0001 O 0 0 0 0 0 0 0 0 0 0 0 0
 0.6072 0.0000 -0.0004 C 0 0 0 0 0 0 0 0 0 0 0 0
 1.1472 0.9353 0.0016 H 0 0 0 0 0 0 0 0 0 0 0 0
 1.1472 -0.9353 0.0016 H 0 0 0 0 0 0 0 0 0 0 0 0
 1 2 2 0 0 0 0
 2 3 1 0 0 0 0
 2 4 1 0 0 0 0
```

Compound

GET
POST
PUT
DELETE

Example 2. Retrieve compound as SMILES

```
$ curl -H "Accept:chemical/x-daylight-smiles"
http://apps.ideaconsult.net:8080/ambit2/compound/1
O=C
```

Example 3. Query compounds

```
$ curl -H Accept:chemical-mime "
http://apps.ideaconsult.net:8080/ambit2/query/compound/{any-identifier-or-keyword}
```

```
$ curl -H Accept:chemical-mime "
http://apps.ideaconsult.net:8080/ambit2/query/smarts?search={smarts}
```

Resources: Dataset

Dataset

Provides access to chemical compounds and their features (e.g. structural, physical-chemical, biological, toxicological properties)

```
@prefix ad: <http://apps.ideaconsult.net:8080/ambit2/dataset/> .
@prefix af: <http://apps.ideaconsult.net:8080/ambit2/feature/> .
@prefix ot: <http://www.opentox.org/api/1.1#> .
...
ad:9 a      ot:Dataset ;
      ot:dataEntry
        [ a      ot:DataEntry ;
          ot:compound
            <http://apps.ideaconsult.net:8080/ambit2/compound/413/conformer/409421> ;
          ot:values
            [ a      ot:FeatureValue ;
              ot:feature af:21576 ;
              ot:value "3.309999942779541"^^xsd:double
            ] ;
          ot:values
            [ a      ot:FeatureValue ;
              ot:feature af:21573 ;
              ot:value "3.0"^^xsd:double
            ]
        ] ;
```

Operations

- POST – Upload a dataset
- PUT – Update the dataset content
- DELETE – Remove the dataset

Representation

RDF/XML (mandatory), MOL, SDF, CSV, TXT, ARFF, .. (optional)

- The dataset consists of data entries.
- Each entry is associated with exactly **one chemical compound**, identified by its URI and available via OpenTox Compound service API;
- One and the same compound can be associated with multiple dataset entries;
- Every “column” is associated with a **Feature**, its representation should be available via OpenTox Feature API

Fea

G

P

P

D

Dataset

GET

POST

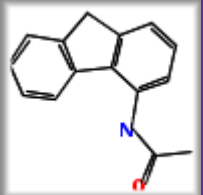
PUT

DELETE

OpenTox datasets: Uniform access to data

Everything described by W3C RDF (Resource Description framework)

Compound/ Data	http://myhost.com/feature/21580	http://myhost.com/feature/21589	http://myhost.com/feature/21573	http://myhost.com/feature/21576	http://myhost.com/feature/21588	http://myhost.com/feature/21858	http://myhost.com/feature/22114
http://myhost.com/compound/413	N,N-dimethyl-4-aminoazobenzene	<chem>CN(C1=CC=C(C=C1)N=N/C2=CC=CC=C2)C</chem>	3	3.31	225.3	YES	3.123
http://myhost.com/compound/44497	4-acetamidofluorene	<chem>CC(=O)Nc1ccc2c(c1)ccc3ccccc23</chem>					
...					



<http://myhost.com/feature/21573>

<http://myhost.com/feature/21858>

<http://myhost.com/feature/22114>

a `ot:Feature` , `ot:NumericFeature` ;

`dc:creator`

"<http://www.blueobelisk.org/ontologies/chemoinformatics-algorithms/#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

Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://apps.ideaconsult.net:8080/ToxPredict/user/admin/report/Stats?header=TRUE

Disable Cookies CSS Forms Images Information Miscellaneous Outline Resize Tools View Source Options

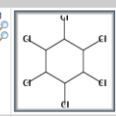
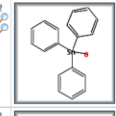
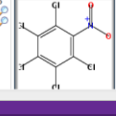
http://apps.ideaconsult.net:8080/ToxPredict/user/admin/report/Stats?header=TRUE

Number of compounds	1. pre-registered substances_20090327.xml	2. CPDBAS: Carcinogenic Potency Database - All Species	3. DBPCAN: EPA Water Disinfection By-Products with Carcinogenicity Estimates	4. ToxCast_ToxRefDB_20091214.txt	5. EPAFHM: EPA Fathead Minnow Acute Toxicity	6. KIERBL: EPA Estrogen Receptor Ki Binding Study (Laws et al.)	7. IRIS: EPA Integrated Risk System (IRIS) Toxicity Review Data	8. FDAMDD: FDA Maximum (Recommended) Daily Dose	9. Burci mutagenicity dataset.sdf	10. c049884m_caco2-training_set.sdf	11. ECETOX Technical Report No. 66 Skin Irritation and corrosion Reference Chemicals data base (1995)	12. ISSMIC v2a_151_2Apr09.sdf	13. Compilation of historical local lymph node assay data for the evaluation of skin
1. pre-registered substances_20090327.xml	143835	259	69	41	33	171	51						
2. CPDBAS: Carcinogenic Potency Database Summary Tables - All Species	259	1515	11	5	59	21	24						
3. DBPCAN: EPA Water Disinfection By-Products with Carcinogenicity Estimates	41	11	108	0	9	1	9						
4. ToxCast_ToxRefDB_20091214.txt	33	59	0	0	307	25	25						
5. EPAFHM: EPA Fathead Minnow Acute Toxicity	171	97	1	13	25	616	18						
6. KIERBL: EPA Estrogen Receptor Ki Binding Study (Laws et al.)	51	34	0	6	25	18	278						
7. IRIS: EPA Integrated Risk System (IRIS) Toxicity Review Data	198	210	2	9	126	93	26						
8. FDAMDD: FDA Maximum (Recommended) Daily Dose	53	150	0	0	1	16	6						
9. Burci mutagenicity dataset.sdf	1740	503	52	36	65	180	57						
10. c049884m_caco2-training_set.sdf	22	23	0	0	0	3	1						
11. ECETOX Technical Report No. 66 Skin Irritation and corrosion Reference Chemicals data base (1995)	138	6	1	1	0	10	0						
12. ISSMIC v2a_151_2Apr09.sdf	136	24	1	0	0	5	1						
13. Compilation of historical local lymph node assay data for the evaluation of skin	170	17	2	1	0	9	1						

Find: issca Next Previous Highlight all Match case Phrase not found

Done

Search results Dataset=112.Dataset= Download as Max number of hits: 100

Compound	ToxCast_To	ToxCast_To	ToxCast_To	ToxCast_To
	CHR_Mouse_Ureter_2_PrenoplasticlesionCHR_Mouse_Nose_1_AnylesionCHR_Rat_Trachea_3_NeoplasticlesionCHR_Mouse_Pit			
1 	1000000.0	1000000.0	1000000.0	1000000.0
2 	1000000.0	1000000.0	1000000.0	1000000.0
3 	NA	NA	1000000.0	NA

Example: mutagenicity dataset

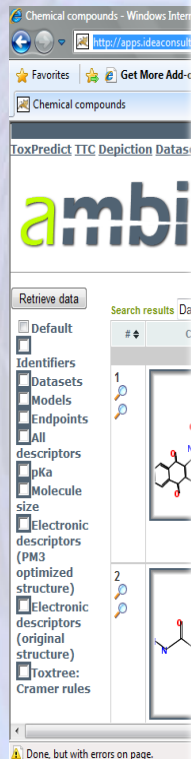
- <http://apps.ideaconsult.net:8080/ambit2/dataset/2344> (the dataset)
- <http://apps.ideaconsult.net:8080/ambit2/dataset/2344/metadata> (metadata, obviously)

The screenshot displays the Ambit web application interface within a Windows Internet Explorer browser. The address bar shows the URL <http://apps.ideaconsult.net:8080/ambit2/dataset/2344/metadata>. The page features a navigation menu with links such as ToxPredict, TTC, Depiction, Datasets, Chemical compounds, Similarity, Substructure, Algorithms, References, Features, Templates, Models, Ontology, RDF, playground, and Help. The main content area includes a search bar with fields for SMILES and Keywords, and a "Search" button. Below the search bar, a message states: "Search for substructure and properties. This site and AMBIT REST services are under development!". The "Retrieve data" section shows "Search results Dataset = 2344" and a "Download as" button. A table of results is displayed with columns for Compound, tox_bench, MC, Example, WDI Name, canonical smiles, Source, and to. The table contains two entries, each with a chemical structure image, a toxicity score of 0.0, and a source reference.

#	Compound	tox_bench	MC	Example	WDI Name	canonical smiles	Source	to
1		0.0	+			<chem>O=C1C2CCCC2C(=O)C3CCCC4C3N1C5C6C(=O)C7CCCC7C(=O)C6C8N1C9C%10C(=O)C%11CCCC%11C(=O)C%10CCC9C8C45</chem>	VITIC	JUDSON, P.N., DOERRER, N.G., HANZLIK, R.P., HARTMANN, J., HOLDER, J. M., HARTMANN, J., SMITH, M., TH AND ZEIGER, I. CREATION OF TOXICOLOGY CENTRE, TOX 2):117-28, 20
2		0.0	+			<chem>NNC(=O)CNC(=O)C=N#N</chem>	CCRIS	MCCANN, J. C AND AMES, B. CARCINOGEN THE SALMON TEST: ASSAY PROC. NATL. (12):5135-51

Example: mutagenicity dataset

Activity



@prefix ot: <http://www.opentox.org/api/1.1#> .
@prefix dc: <http://purl.org/dc/elements/1.1/> .
@prefix : <http://apps.ideaconsult.net:8080/ambit2/> .
@prefix ota: <http://www.opentox.org/algorithmTypes.owl#> .
@prefix otee: <http://www.opentox.org/echaEndpoints.owl#> .
...
@prefix af: <http://apps.ideaconsult.net:8080/ambit2/feature/> .

af:28958

a ot:Feature , ot:NumericFeature ;
dc:creator "194.141.0.136" ;
dc:title "Activity" ;
ot:hasSource "tox_benchmark_N6512.sdf" ;
ot:units "" ;
= otee:Mutagenicity .

ot:hasSource
a owl:ObjectProperty .

ot:units
a owl:DatatypeProperty .

ot:Feature
a owl:Class .

ot:NumericFeature
a owl:Class ;
rdfs:subClassOf ot:Feature .

Query: Is there other mutagenicity data available?

<http://apps.ideaconsult.net:8080/ambit2/feature?sameas=http%3A%2F%2Fwww.opentox.org%2FechaEndpoints.owl%23Mutagenicity>

<http://apps.ideaconsult.net:8080/ambit2/feature/21590>

<http://apps.ideaconsult.net:8080/ambit2/feature/21611>

<http://apps.ideaconsult.net:8080/ambit2/feature/26221>

<http://apps.ideaconsult.net:8080/ambit2/feature/28958>

Find	Name	URL	Values
	SAL	ISSCAN_v3a_1153_19Sept08.1222179139.sdf	NO
	ActivityOutcome_CPDBAS_Mutagenicity	CPDBAS_v5d_1547_20Nov2008.sdf	YES
	ActivityOutcome_CPDBAS_Mutagenicity	ambit2_75288.sdf	NO
	ActivityOutcome_CPDBAS_Mutagenicity		
	ActivityOutcome_CPDBAS_Mutagenicity		
	Mutagenic Activity in TA100 (3=active; 1=inactive)	qsar6train.csv	YES
	Ames test categorisation	Burci_mutagenicity_dataset.sdf	YES
	Mutagenic Activity in TA100 (3=active; 1=inactive)	270.sdf	NO
	Mutagenic Activity in TA100 (3=active; 1=inactive)		
	Prediction feature for http://apps.ideaconsult.net:8080/ambit2/feature/26701 endpoint prediction	/OpenTox-dev/model/10mOpenToxModel_j48_30	
	Activity	tox_benchmark_N6512.sdf	NO

Merge mutagenicity data

- `http://apps.ideaconsult.net:8080/ambit2/dataset/2344?feature_uris[]=http://apps.ideaconsult.net:8080/ambit2/feature/28958&feature_uris[]=http://apps.ideaconsult.net:8080/ambit2/feature/21611&feature_uris[]=http://apps.ideaconsult.net:8080/ambit2/feature/26221&feature_uris[]=http://apps.ideaconsult.net:8080/ambit2/feature/21590`

The screenshot shows the Ambit web application interface in a Mozilla Firefox browser. The address bar displays a URL with multiple feature URIs. The page header includes navigation links like 'oxPredict', 'TTC', 'Depiction', 'Datasets', 'Chemical compounds', 'Similarity', 'Substructure', 'Algorithms', 'References', 'Features', 'Templates', 'Models', 'Ontology', 'RDF', 'playground', and 'Help'. The main content area features the 'ambit' logo, search input fields for 'SMARTS' and 'Keywords', and a 'Search' button. Below the search area, a table displays search results for dataset 2344. The table has columns for 'Compound', 'tox_benchm', 'CPDBAS v5d', and 'Burci muta'. The first two columns are further divided into 'Activity' and 'ActivityOutcome CPDBAS Mutagenicity'. The table lists two compounds: one with a complex polycyclic structure and an activity of 0.0, and another with a linear structure and an activity of 1.0, which is categorized as 'mutagen'.

#	Compound	tox_benchm	CPDBAS v5d	Burci muta
		Activity	ActivityOutcome CPDBAS Mutagenicity	Ames test categorisation
1		0.0		
2		1.0		mutagen

Dataset : metadata and features

Description	URI Template
Retrieve entire dataset content. If uri-list, retrieve only compound URIs	http://host:port/dataset/{id}
Retrieve representation of features (columns) of the dataset	http://host:port/dataset/{id}/feature
Retrieves dataset metadata (name, etc.)	http://host:port/dataset/{id}/metadata

```
$ curl -H "Accept:application/rdf+xml" http://apps.ideaconsult.net:8080/ambit2/dataset/9/metadata
<rdf:RDF
xmlns:ot="http://www.opentox.org/api/1.1#"
.....
  xmlns:rdfs="http://www.w3.org/2000/01/rdf-schema#"
  xml:base="http://apps.ideaconsult.net:8080/ambit2/">
<ot:Dataset rdf:about="dataset/9">
  <dc:source>ISSCAN_v3a_1153_19Sept08.1222179139.sdf</dc:source>
  <dc:publisher>somebody</dc:publisher>
  <rdfs:seeAlso>
    <bx:Entry rdf:about="reference/20117">
      <rdfs:seeAlso>http://www.epa.gov/NCCT/dsstox/sdf_isscan_external.html</rdfs:seeAlso>
      <dc:title>ISSCAN_v3a_1153_19Sept08.1222179139.sdf</dc:title>
    </bx:Entry>
  </rdfs:seeAlso>
  <dc:title>ISSCAN: Istituto Superiore di Sanita, CHEMICAL CARCINOGENS: STRUCTURES AND EXPERIMENTAL DATA</dc:title>
</ot:Dataset>
</rdf:RDF>
```

Data publishing

1) POST a file with chemical structures and properties to OpenTox dataset service.

- The structures and data are assigned a dataset URL and become available by multiple formats (RDF, Chemical MIME, CSV, Weka ARFF)

2) Assign metadata

- PUT /dataset/{id}/metadata

3) Annotate any of dataset features /dataset/{id}/feature by assigning links to relevant ontologies

- PUT /feature/{id}

HTTP GET / POST

Find chemical compounds,
return dataset URL;
<http://host2/compound/id>

Upload data, receive
dataset URL
<http://host2/dataset/id>

Dataset service

Structures,
endpoints
& predictions

Annotation

Ontology service

Published models,
Algorithms,
Ontologies,
metadata

Algorithms
ontologies

Toxicology related
ontologies

Dataset and Ontology - find an assay, linked to specific gene

PREFIX ot:<<http://www.opentox.org/api/1.1#>>
 PREFIX ota:<<http://www.opentox.org/algorithms.owl#>>
 PREFIX owl:<<http://www.w3.org/2002/07/owl#>>
 PREFIX dc:<<http://purl.org/dc/elements/1.1/>>
 PREFIX rdfs:<<http://www.w3.org/2000/01/rdf-schema#>>
 PREFIX rdf:<<http://www.w3.org/1999/02/22-rdf-syntax-ns#>>
 PREFIX otee:<<http://www.opentox.org/echaEndpoints.owl#>>
 PREFIX toxcast:<<http://www.opentox.org/toxcast#>>

select ?Feature ?title ?id ?assay ?geneid ?gene

where {

 ?Feature rdf:type ot:Feature.

 {?Feature dc:title ?title}.

 {?Feature owl:sameAs ?assay}.

 {?assay toxcast:gene ?geneid}.

 {?assay toxcast:hasProperty ?genename}.

 {?genename rdf:type toxcast:GENE_NAME}.

}



Query an OpenTox ontology service at
<http://ambit.uni-plovdiv.bg:8082/ontology>

Chemical compounds

Search results Dataset = 961 hits: 100

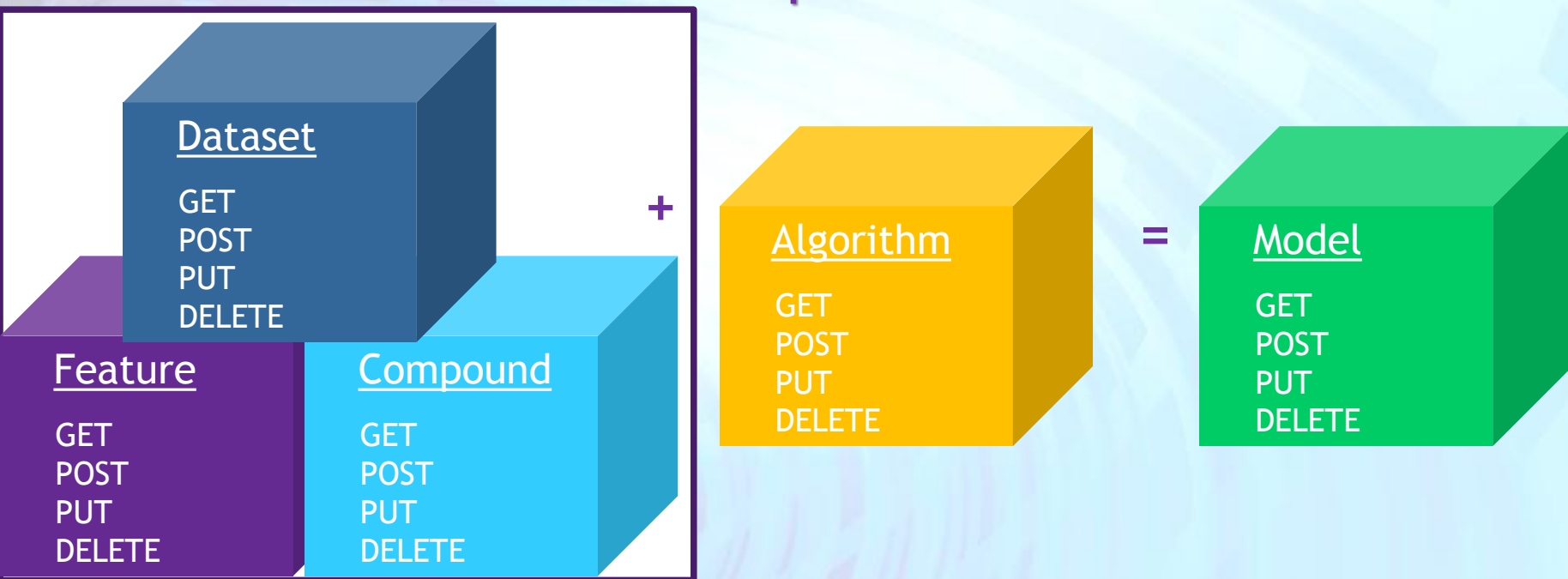
#	Compound	ToxCast At	Benigni /	Benigni /
		ATG RORE CIS	Structural Alert for genotoxic carcinogenicity	Structural Alert for nongenotoxic carcinogenicity
1		1000000.0	NO	NO
2		1000000.0	NO	NO
3		1000000.0	NO	NO
4		1000000.0	NO	NO

?feat

126

OpenTox dataset : create a model

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



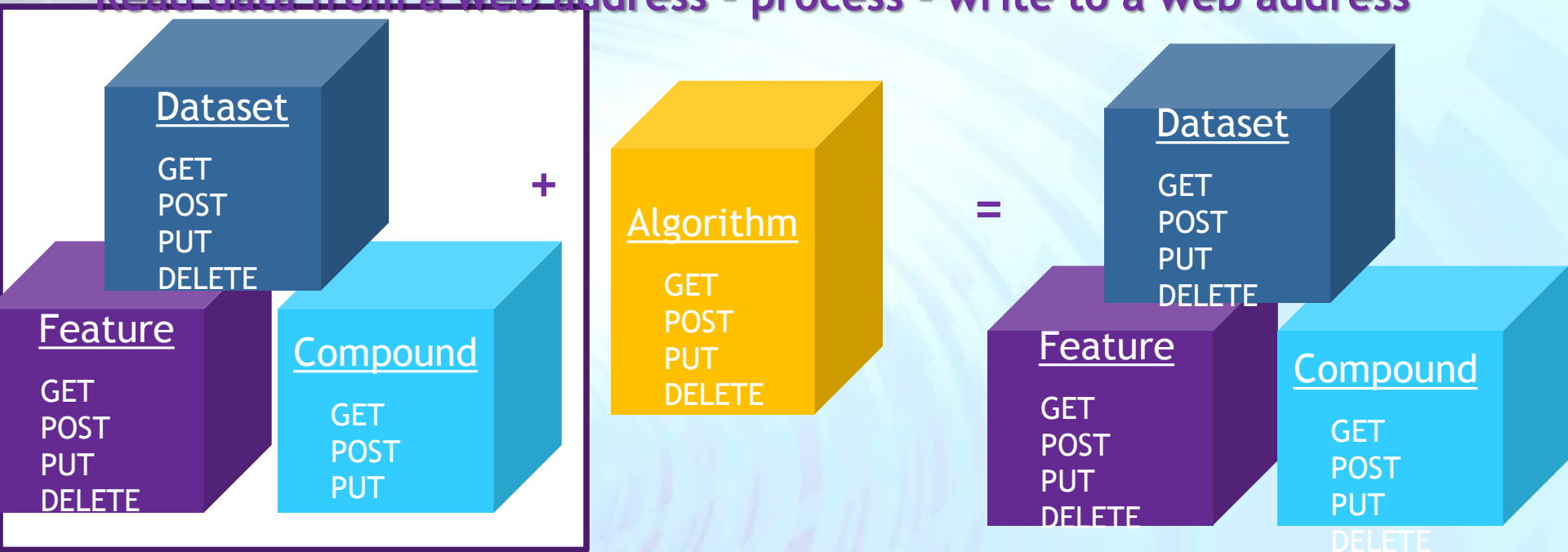
<http://myhost.com/algorithm/neuralnetwork>

<http://myhost.com/dataset/trainingset1>

<http://myhost.com/model/predictivemodel1>

OpenTox dataset : descriptor calculation

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



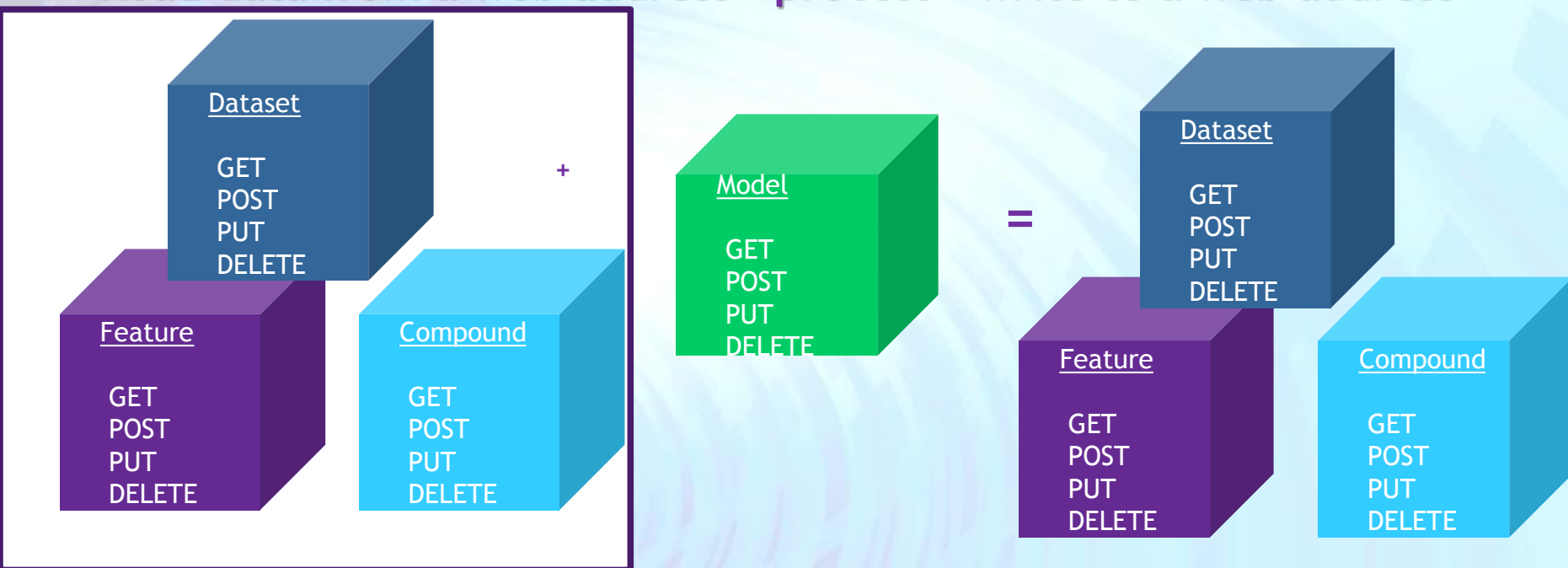
<http://myhost.com/algorithm/{descriptorX}>

<http://myhost.com/dataset/trainingset1>

<http://myhost.com/dataset/results>

OpenTox dataset: apply a model

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



<http://myhost.com/model/predictivemodel1>

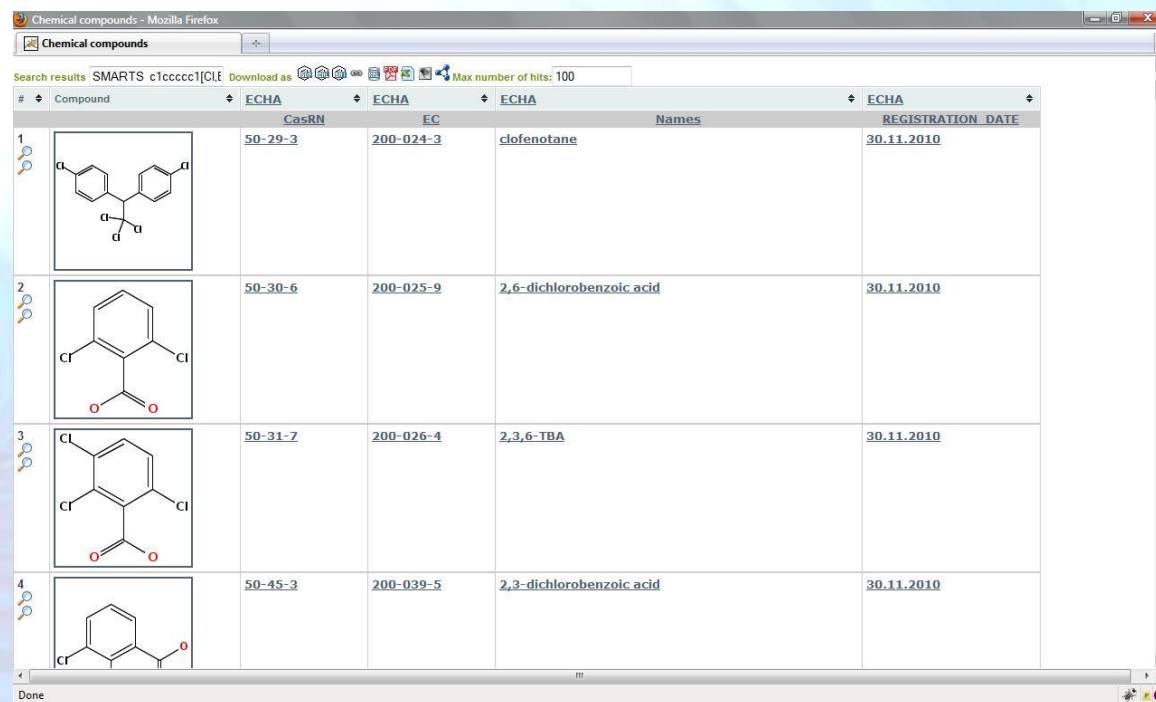
<http://myhost.com/dataset/id1>

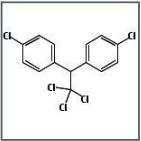
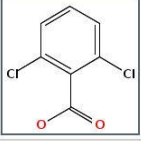
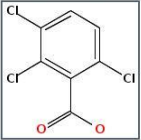
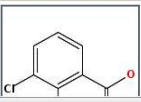
<http://myhost.com/dataset/results1>

OpenTox datasets : Substructure and similarity search

REST web service interface

[http://apps.ideaconsult.net:8080/ambit2/query/smarts?search=c1ccccc1\[Cl,Br,F,I\]](http://apps.ideaconsult.net:8080/ambit2/query/smarts?search=c1ccccc1[Cl,Br,F,I])



Compound	ECHA	EC	Names	REGISTRATION DATE
	50-29-3	200-024-3	clofenotane	30.11.2010
	50-30-6	200-025-9	2,6-dichlorobenzoic acid	30.11.2010
	50-31-7	200-026-4	2,3,6-TBA	30.11.2010
	50-45-3	200-039-5	2,3-dichlorobenzoic acid	30.11.2010

Datasets : Structure and quality labels

Dataset	OK	Probably OK	Probably ERROR	Unknown	Probably ERROR%
ECHA list of pre-registered substances	N/A	N/A	N/A	N/A	N/A
Chemical Identifier Resolver	67779	5314	3638	3471	4.75%
ChemIDplus	64802	7986	921	1745	1.24%
ChemDraw	17918	1147	502	478	2.57%
JRC PRS list	61332	4833	4022	2880	5.83%
ISSCAN	931	50	98	62	9.40%
CPDBAS	778	37	0	693	0%
DBPCAN	60	2	0	147	0%
EPAFHM	281	5	0	331	0%
KIERBL	102	1	0	175	0%
IRISTR	346	16	0	177	0%
FDAMDD	213	19	1	983	0.08%
ECETOC skin irritation	158	12	0	5	0%
Skin sensitisation (LLNA)	160	7	4	38	1.95%
Bioconcentration factor (BCF) Gold Standard Database	N/A	N/A	N/A	N/A	N/A

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

Dataset
Resource

Descriptor
resource

Assay
resource

Chemical
compound

Blue Obelisk
algorithms
ontology

Regression
Classification
Quantum
Chemistry
Descriptors, etc.

OpenTox
algorithm types
ontology

Toxicology related
ontologies

http://apps.ideaconsult.net:8080/ambit2/dataset/R545

data Entry	compound
values	http://apps.ideaconsult.net:8080/ambit2/compound/38/conformer/419609
feature	TopoPSA
Feature value	6.480000019073486
type	Feature Value
feature	nHBDon
Feature value	0.0
type	Feature Value
feature	caco2
Feature value	8.849999904632568
type	Feature Value
feature	WNSA-3
Feature value	-374800205230713
type	Feature Value
feature	FP5A-2
Feature value	0.8797000050544739
type	Feature Value
compound	http://apps.ideaconsult.net:8080/ambit2/compound/144824/conformer/419615
values	TopoPSA
Feature value	210.5399938861328
type	Feature Value
feature	nHBDon
Feature value	5.0
type	Feature Value
feature	caco2
Feature value	-5.920000076293945
type	Feature Value
feature	WNSA-3
Feature value	-64.12879943847656
type	Feature Value
feature	FP5A-2
Feature value	2.147799968719482
type	Feature Value

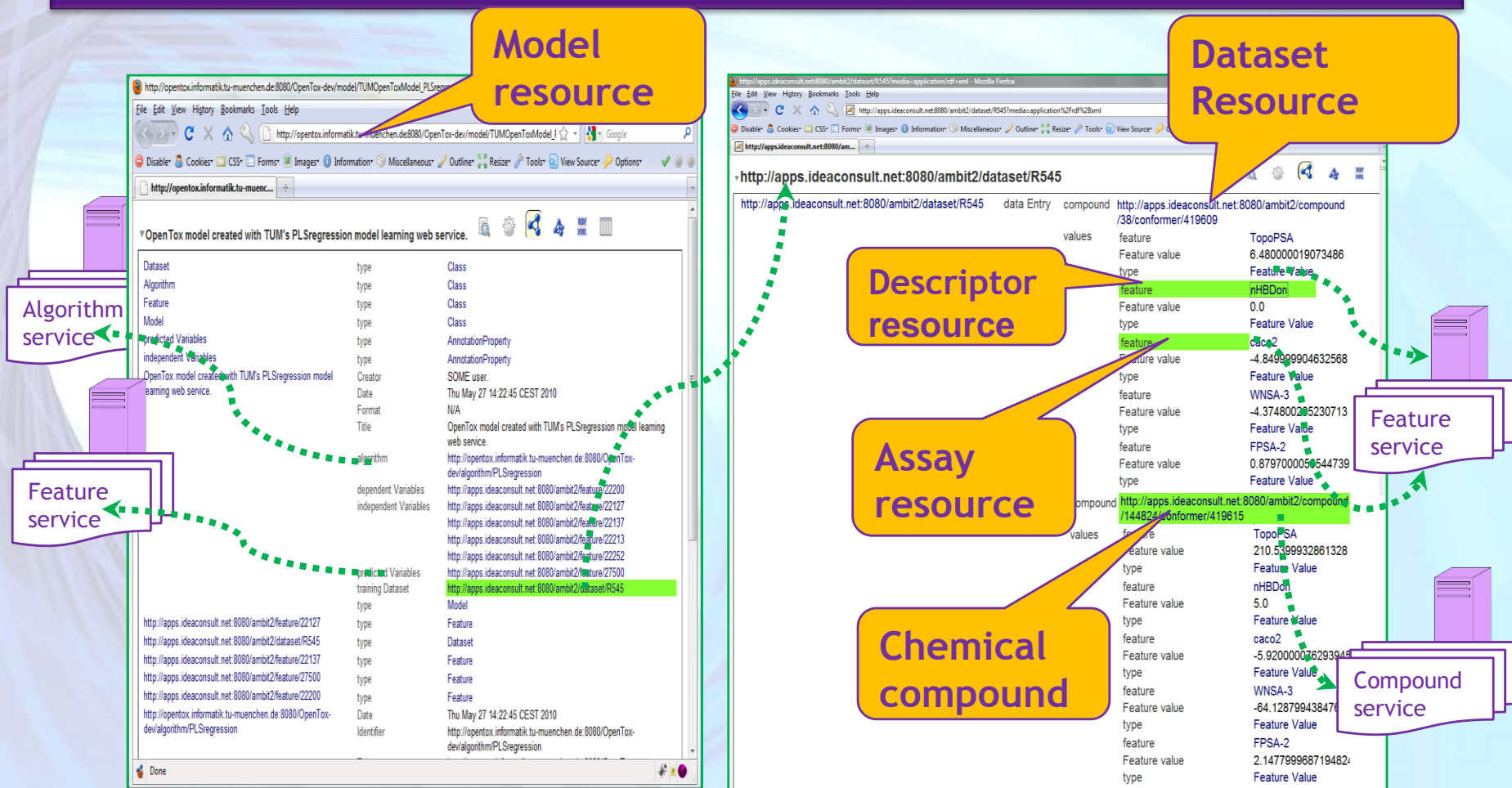
http://apps.ideaconsult.net:8080/ambit2/feature/22213

Name of the algorithm	type	Class
type	type	Class
Numeric Feature	type	Class
Source	type	subClassOf
Units	type	ObjectProperty
nHBDon	type	DatatypeProperty
Source	type	http://www.blueobelisk.org/ontologies/chemoinformatics-algorithms/#nHBDon
Units	type	nHBDon
Source	type	http://apps.ideaconsult.net:8080/ambit2/algorithm/org.openscience.cdk.qsar.descriptors.molecular.HBondDonorCountDescr
Units	type	Numeric Feature

http://apps.ideaconsult.net:8080/ambit2/feature/22200

Numeric Feature	type	Class
type	type	Class
Source	type	subClassOf
Units	type	ObjectProperty
caco2	type	DatatypeProperty
Source	type	caco2
Units	type	c049084m_caco2-training_set.sdf
Source	type	Numeric Feature
Units	type	Gastrointestinal absorption

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



Development of Predictive Toxicology Applications

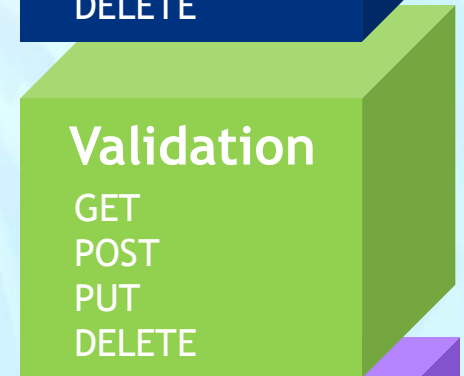
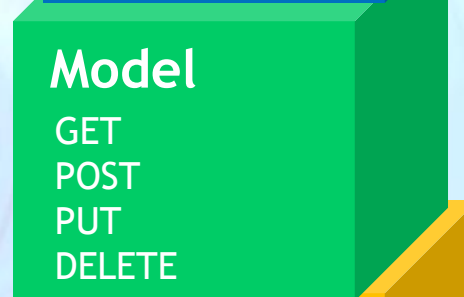
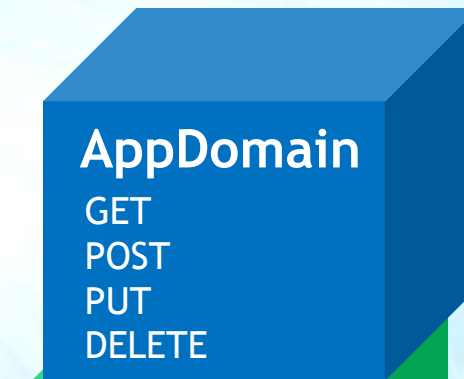
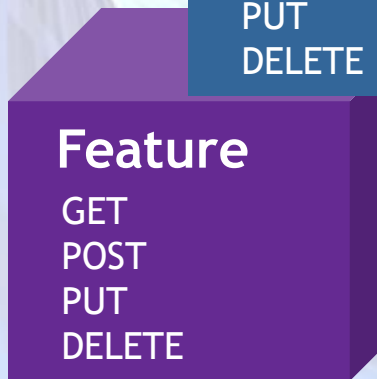
An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Algorithms

Stefan Kramer

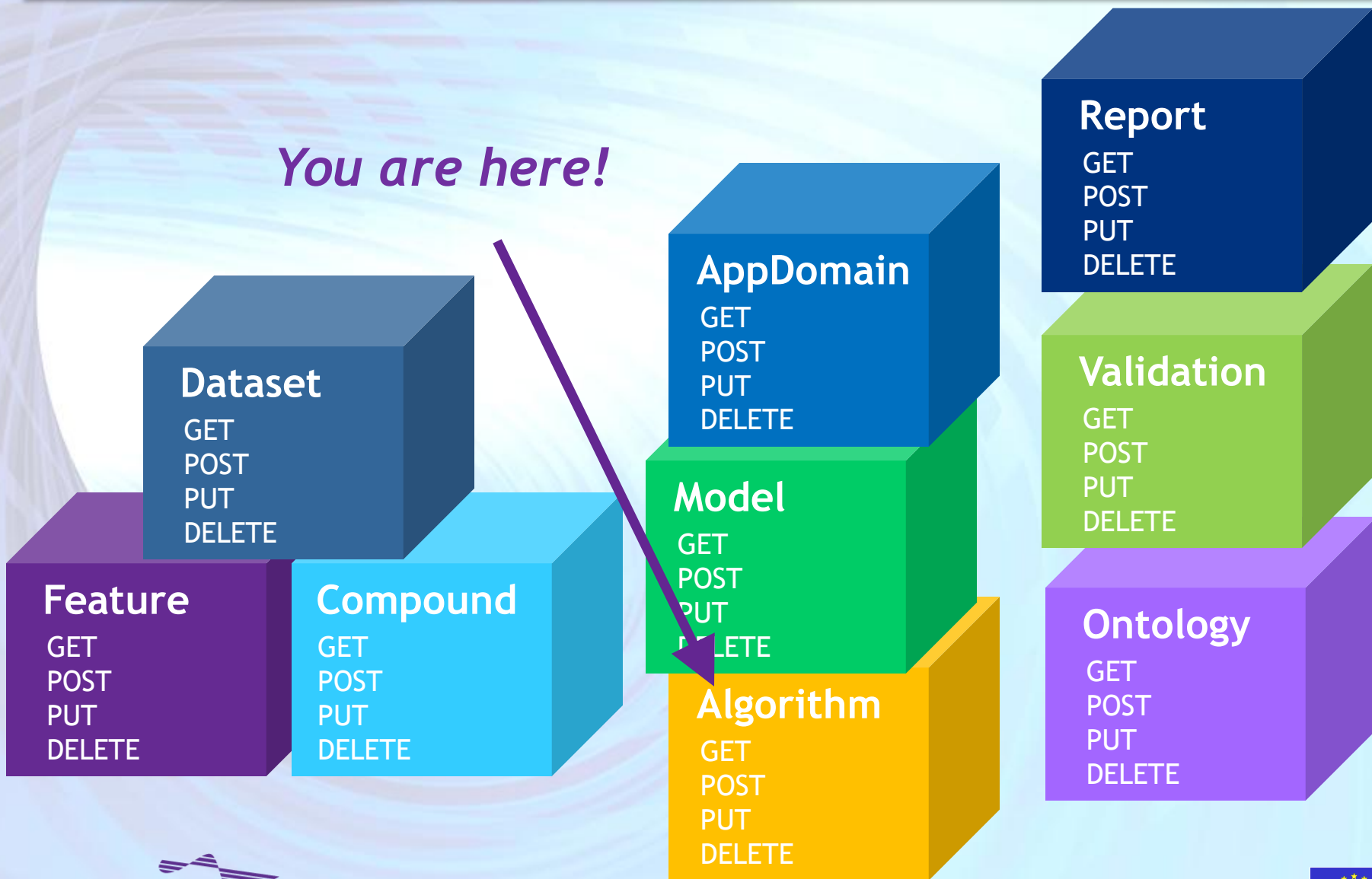
(Technische Universität München, Munich, Germany)

Overview of Application Programming Interfaces



Overview of Application Programming Interfaces

You are here!



Overview of Algorithms in the OpenTox Framework

- Algorithms for **descriptor calculation: generation and selection of features** for the representation of chemicals (structure based features, chemical and biological properties),
- **Classification and regression** algorithms for the creation of (Q)SAR models,
- 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 domains, read across,)

OpenTox Algorithms: Descriptor Calculation and Feature Selection

- **Descriptor calculation:** services based on
 - OpenBabel
 - JoeLib2
 - CDK
 - multi-level neighborhood of atoms (MNA)
 - substructure/fragment generation (“product line”: gSpan, FreeTreeMiner, BBRCs, LastPM; details later)
- **Feature selection:**
 - service 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 selection, ...

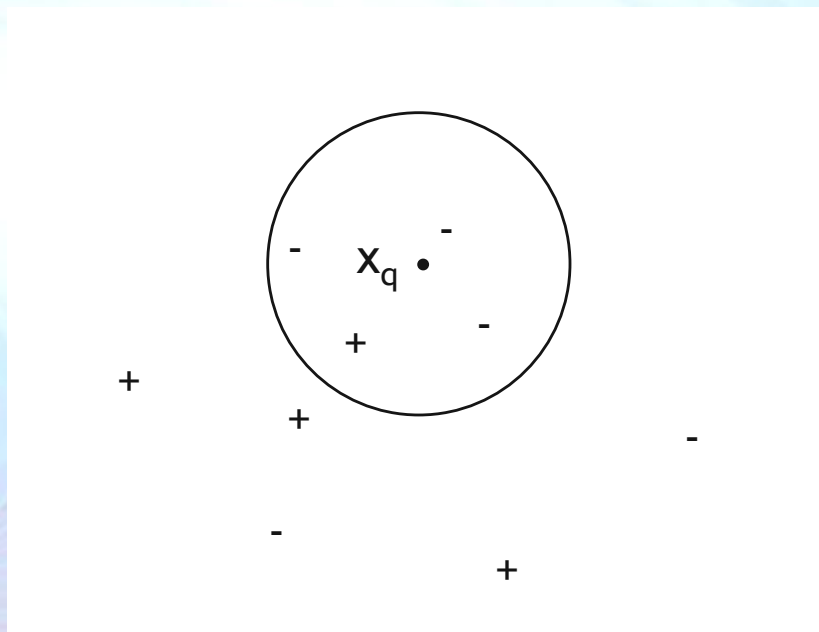
OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**

*Leaving out algorithms/approaches for
applicability domain and validation for now*

OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**



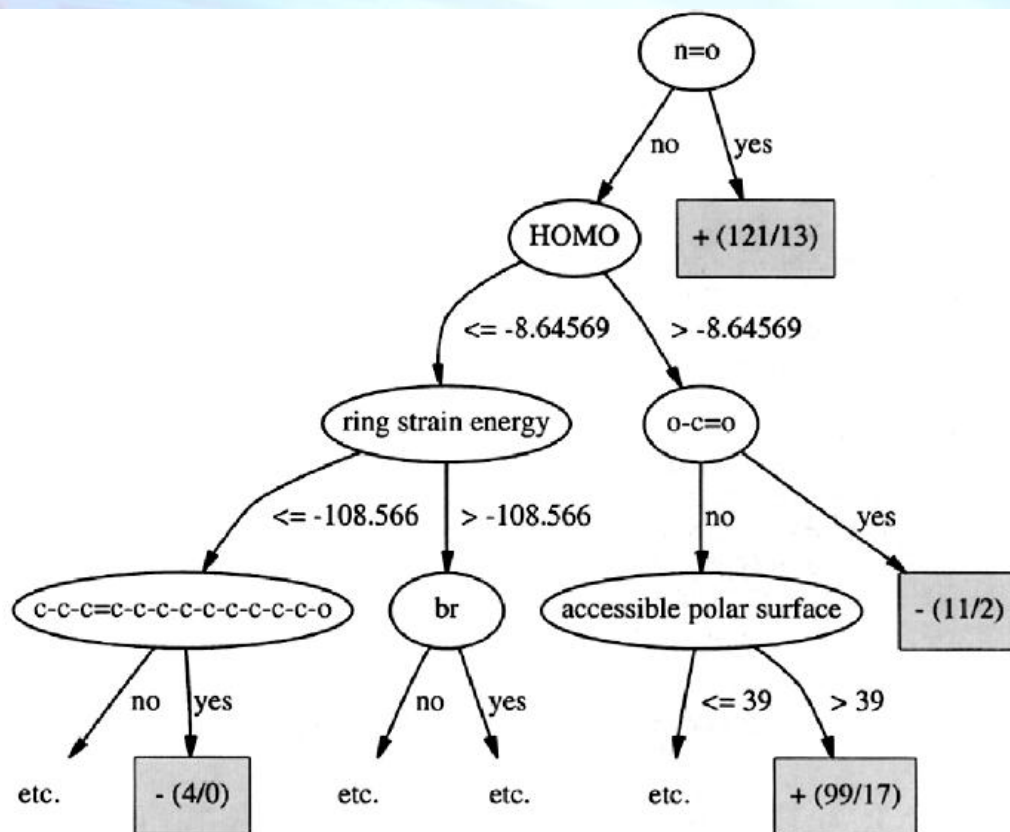
Leaving out algorithms/approaches for applicability domain and validation for now

OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning algorithms:
 - **decision trees (J48)**

OpenTox Algorithms: Classification / SAR

- Simple
- Machine learning
- decision tree



Leaving out algorithms/approaches for applicability domain and validation for now

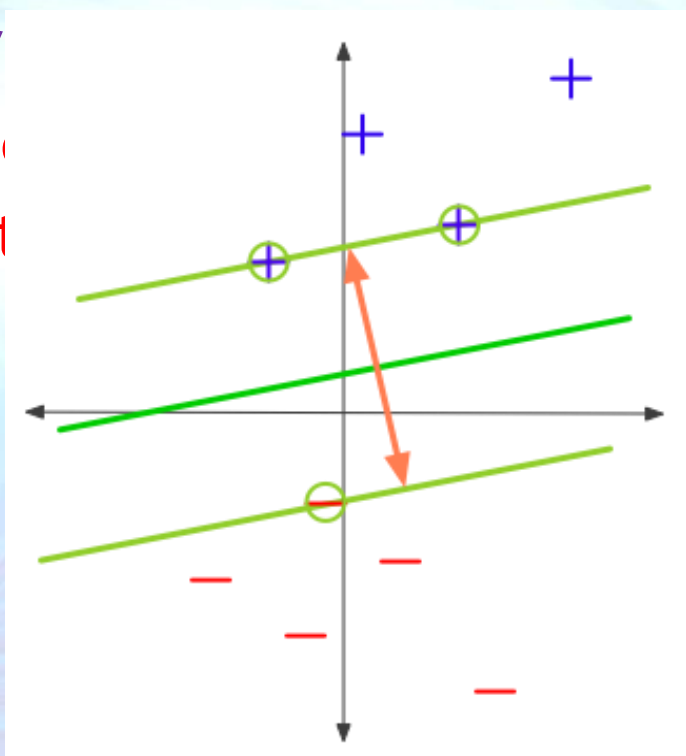
OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning algorithms:
 - **decision trees (J48)**
 - **support vector machines (SVMs)**

*Leaving out algorithms/approaches for
applicability domain and validation for now*

OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning
 - decision tree
 - support vector



*Leaving out algorithms/approaches for
applicability domain and validation for now*

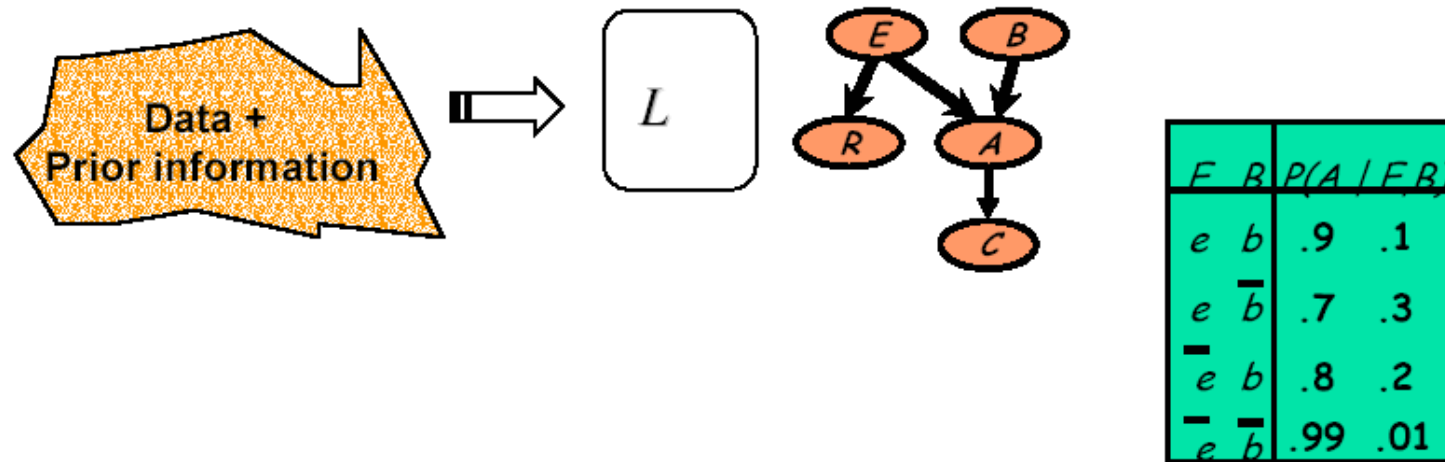
OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning algorithms:
 - **decision trees (J48)**
 - **support vector machines (SVMs)**
- Probabilistic/graphical models
 - **Bayesian network**

Leaving out algorithms/approaches for applicability domain and validation for now

OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning algorithms:
 - **decision trees (J48)**



Leaving out algorithms/approaches for applicability domain and validation for now

OpenTox Algorithms: Classification / SAR

- Simple baseline: **k-Nearest Neighbor**
- Machine learning algorithms:
 - **decision trees (J48)**
 - **support vector machines (SVMs)**
- Probabilistic/graphical models
 - **Bayesian network**

Leaving out algorithms/approaches for applicability domain and validation for now

OpenTox Algorithms: Regression / QSAR

- Simple baseline: **k-Nearest Neighbor**
- Classical statistical algorithms:
 - multiple linear regression (MLR)
 - partial least squares (PLS)
- Machine learning algorithms:
 - **model trees (M5')**

OpenTox Algorithms: Regression / QSAR

- Simple baseline: **k-Nearest Neighbor**
- Classical statistical algorithms:
 - **multiple linear regression (MLR)**

— `log_fluence <= -6.01 :`
| `log_hr321 <= -0.112 : LM1`
| `log_hr321 > -0.112 : LM2`
`log_fluence > -6.01 :`
— | `log_hr321 <= 0.0846 : LM3`
| `log_hr321 > 0.0846 : LM4`

LM1: $\log_{t90} = -0.879 + 0.0353\log_{hr321} - 0.373\log_{fluence}$
+ $0.0394mfbmfr_class=inter, long + 0.0327mfbmfr_class=long$
LM2: $\log_{t90} = 0.00965 - 0.138\log_{hr321} - 0.203\log_{fluence}$
+ $0.0394mfbmfr_class=inter, long + 0.0327mfbmfr_class=long$

OpenTox Algorithms: Regression / QSAR

- Simple baseline: **k-Nearest Neighbor**
- Classical statistical algorithms:
 - **multiple linear regression (MLR)**
 - **partial least squares (PLS)**
- Machine learning algorithms:
 - **model trees (M5')**
 - **support vector regression**
- Probabilistic/graphical models:
 - **Gaussian process regression**

OpenTox API for Algorithms

Description	Method	URI	Parameters	Result	Status codes
Get URIs of all available algorithms	GET	/algorithm	(optional) ?sameas=URI-of-the-owl:sameAs-entry	List of all algorithm URIs or RDF representation, or algorithms of specific types, if query parameter exists. Returns all algorithms, for which owl:sameAs is given by the query.	200,404,503
Get the ontology representation of an algorithm	GET	/algorithm/{id}	–	Algorithm representation in one of the supported MIME types.	200,404,503
Apply the algorithm	POST	/algorithm/{id}	dataset_uri parameter prediction_feature , more to be specified and documented by algorithm provider dataset_service =data setserviceuri	<i>model URI</i> <i>dataset URI</i> <i>featureURI</i> Redirect to task URI for time consuming computations.	200,303,404,503

GET <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/>

TUM OpenTox REST web services - Mozilla Firefox

Datei Bearbeiten Ansicht Chronik Lesezeichen Extras Hilfe

<http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/>

Meistbesuchte Seiten Erste Schritte Aktuelle Nachrichten news.ORF.at sport.ORF.at SPIEGEL ONLINE - Nac... derStandard.at sueddeutsche.de Topt... SPIEGEL ONLINE - Uni...

TUM OpenTox REST web services

TUM Technische Universität München **OpenTox**

TUM - OpenTox - REST services 1.1

This site and [TUM](http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/) OpenTox REST web services are under development!
The full [API](#) can be found on the opentox.org website

Available algorithms:

- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/kNNclassification>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/J48>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/kNNregression>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/PLRegression>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/M5P>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/GaussP>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/LR>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/BayesNet>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/FTM>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/gSpan>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/FTM/{smiles}>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/gSpan/{smiles}>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/CDKPhysChem>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/JOELIB2>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/InfoGainAttributeEval>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/PrincipalComponents>
- <http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/ChiSquared>

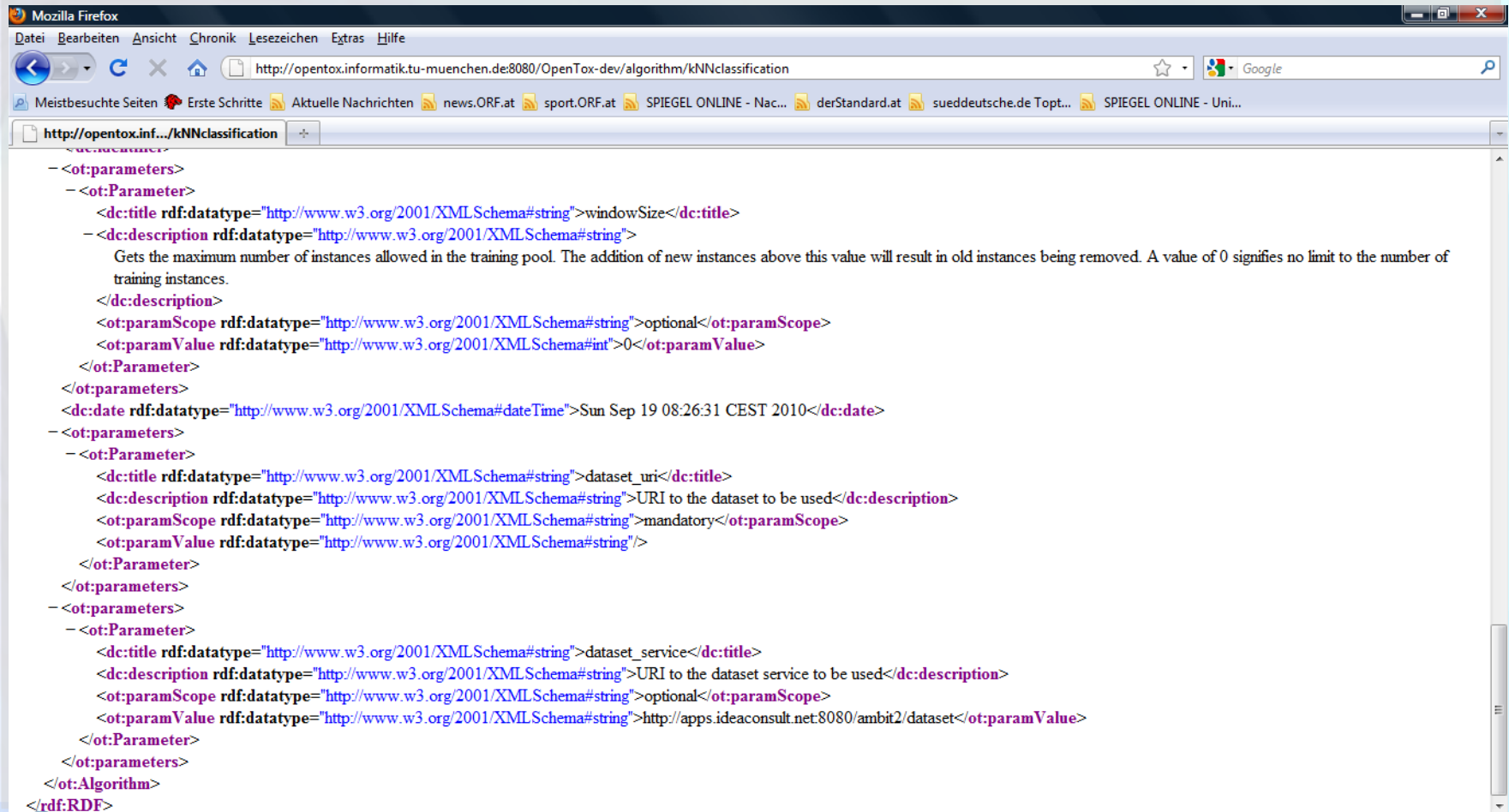
[Initial documentation](#)

k-NN Classification

```
Mozilla Firefox
Datei Bearbeiten Ansicht Chronik Lesezeichen Extras Hilfe
http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/kNNclassification
Meistbesuchte Seiten Erste Schritte Aktuelle Nachrichten news.ORF.at sport.ORF.at SPIEGEL ONLINE - Nac... derStandard.at sueddeutsche.de Topt... SPIEGEL ONLINE - Uni...
http://opentox.inf.../kNNclassification

- <rdf:RDF>
  <owl:Class rdf:about="http://www.opentox.org/api/1.1#Algorithm"/>
  <owl:Class rdf:about="http://www.opentox.org/api/1.1#Parameter"/>
  <ot:Algorithm rdf:about="http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/kNNclassification">
    <dc:contributor>joerg.wicker@in.tum.de</dc:contributor>
  <ot:parameters>
    <ot:Parameter>
      <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">nearestNeighbourSearchAlgorithm</dc:title>
      <dc:description rdf:datatype="http://www.w3.org/2001/XMLSchema#string">
        The nearest neighbour search algorithm to use (Default: weka.core.neighboursearch.LinearNNSearch).
      </dc:description>
      <ot:paramScope rdf:datatype="http://www.w3.org/2001/XMLSchema#string">optional</ot:paramScope>
      <ot:paramValue rdf:datatype="http://www.w3.org/2001/XMLSchema#string">LinearNNSearch</ot:paramValue>
    </ot:Parameter>
  </ot:parameters>
  <ot:isA>
    http://www.opentox.org/algorithms.owl#ClassificationLazySingleTarget
  </ot:isA>
  <ot:parameters>
    <ot:Parameter>
      <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">crossValidate</dc:title>
      <dc:description rdf:datatype="http://www.w3.org/2001/XMLSchema#string">
        Whether hold-one-out cross-validation will be used to select the best k value
      </dc:description>
      <ot:paramScope rdf:datatype="http://www.w3.org/2001/XMLSchema#string">optional</ot:paramScope>
      <ot:paramValue rdf:datatype="http://www.w3.org/2001/XMLSchema#boolean">>false</ot:paramValue>
    </ot:Parameter>
  </ot:parameters>
  <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">kNNclassification</dc:title>
  <dc:creator>tobias.eirschick@in.tum.de</dc:creator>
```

k-NN Classification



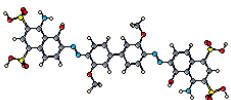
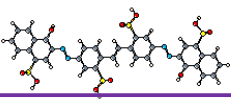
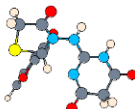
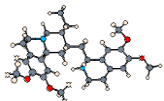
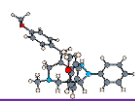
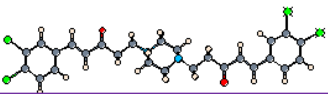
The screenshot shows a Mozilla Firefox browser window with the address bar displaying `http://opentox.informatik.tu-muenchen.de:8080/OpenTox-dev/algorithm/kNNclassification`. The browser's bookmarks bar shows several sites including `news.ORF.at`, `sport.ORF.at`, `SPIEGEL ONLINE`, `derStandard.at`, `sueddeutsche.de`, and `SPIEGEL ONLINE - Uni...`. The main content area displays an RDF document for k-NN classification, which is a subset of the OpenTox ontology. The document is structured as follows:

```
<?xml version="1.0" encoding="UTF-8"?>
<rdf:RDF xmlns:ot="http://www.w3.org/2001/XMLSchema#string" xmlns:dc="http://www.w3.org/2001/XMLSchema#string">
  <ot:parameters>
    <ot:Parameter>
      <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">windowSize</dc:title>
      <dc:description rdf:datatype="http://www.w3.org/2001/XMLSchema#string">
        Gets the maximum number of instances allowed in the training pool. The addition of new instances above this value will result in old instances being removed. A value of 0 signifies no limit to the number of training instances.
      </dc:description>
      <ot:paramScope rdf:datatype="http://www.w3.org/2001/XMLSchema#string">optional</ot:paramScope>
      <ot:paramValue rdf:datatype="http://www.w3.org/2001/XMLSchema#int">0</ot:paramValue>
    </ot:Parameter>
  </ot:parameters>
  <dc:date rdf:datatype="http://www.w3.org/2001/XMLSchema#dateTime">Sun Sep 19 08:26:31 CEST 2010</dc:date>
  <ot:parameters>
    <ot:Parameter>
      <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">dataset_uri</dc:title>
      <dc:description rdf:datatype="http://www.w3.org/2001/XMLSchema#string">URI to the dataset to be used</dc:description>
      <ot:paramScope rdf:datatype="http://www.w3.org/2001/XMLSchema#string">mandatory</ot:paramScope>
      <ot:paramValue rdf:datatype="http://www.w3.org/2001/XMLSchema#string"/>
    </ot:Parameter>
  </ot:parameters>
  <ot:parameters>
    <ot:Parameter>
      <dc:title rdf:datatype="http://www.w3.org/2001/XMLSchema#string">dataset_service</dc:title>
      <dc:description rdf:datatype="http://www.w3.org/2001/XMLSchema#string">URI to the dataset service to be used</dc:description>
      <ot:paramScope rdf:datatype="http://www.w3.org/2001/XMLSchema#string">optional</ot:paramScope>
      <ot:paramValue rdf:datatype="http://www.w3.org/2001/XMLSchema#string">http://apps.ideaconsult.net:8080/ambit2/dataset</ot:paramValue>
    </ot:Parameter>
  </ot:parameters>
</ot:Algorithm>
</rdf:RDF>
```

Development of Novel Algorithms

- Substructure / fragment generation algorithms („product line“)
 - FreeTreeMiner
 - BBRCs (backbone refinement classes)
 - LastPM (latent structure pattern mining)
- Structural clustering and local models
- Fast conditional density estimation for QSAR:
 - quantifying uncertainty in QSAR, confidence intervals, ...

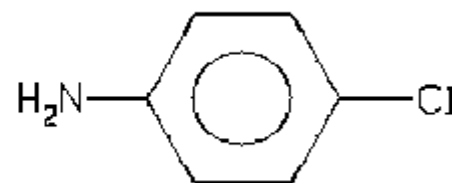
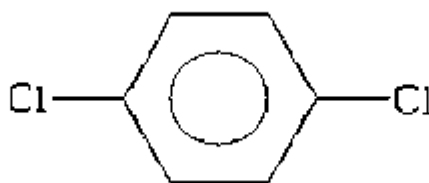
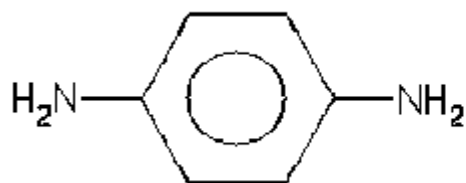
Substructure/Fragment Generation

$C_{A,1}$		CA
$C_{A,2}$		CA
$C_{A,3}$		CA
...
$C_{I,1}$		CI
$C_{I,2}$		CI
$C_{I,3}$		CI
...

CA: confirmed active
CI: confirmed inactive

- First step:
computation of
descriptors
- Computed physico-chemical
properties?
- Predefined
functional groups?
- ...?

Substructure/Fragment Generation



- *Path* patterns
- Minimum frequency = 2
- Just *most specific patterns* used here

$x1 =_{\text{def}} \text{N-c:c:c:c:c:c}$

$x2 =_{\text{def}} \text{Cl-c:c:c:c:c:c}$

$x1$	$x2$	<i>Class</i>
true	false	-
false	true	+
true	true	+

Substructure/Fragment Generation

Statistical learning schemes like SVMs are very good at combining substructures as features into (Q)SAR models

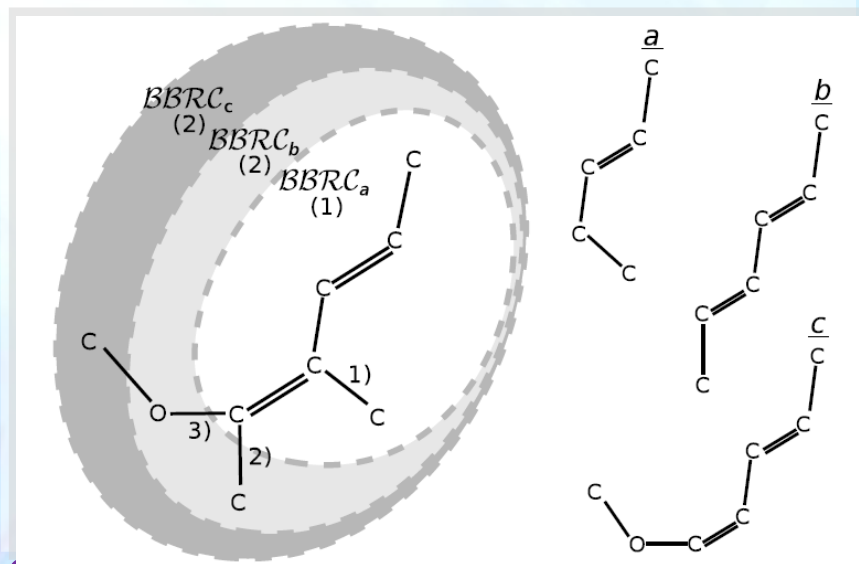
$$f(\mathbf{x}) = \text{sgn}\left(\sum_j^m \gamma_j x_j + c\right) =$$

$\text{sgn}(+1.63 * \text{'c:c:c:c:c:c:c:c:c' (x)}$
 $+1.44 * \text{'C-Cl' (x)}$
 $+1.32 * \text{'C-C-C-C-N-C' (x)}$
 $+1.31 * \text{'C-C-C-O' (x)}$
 $+0.95 * \text{'C-C=C' (x)}$
 $+0.87 * \text{'c:c:c:c:c:n' (x)}$
 $+0.82 * \text{'C-C-C-C=C' (x)}$
 $+0.82 * \text{'C-C-C-N-C' (x)}$
 $+0.80 * \text{'c:c:c-C=O' (x)}$
 $+0.78 * \text{'C-N-C' (x)}$
 $+...$

$-1.48 * \text{'Cl-C-Cl' (x)}$
 $-1.45 * \text{'C-C-C=C-C' (x)}$
 $-1.01 * \text{'C-N-c:c' (x)}$
 $-1.01 * \text{'C-N-c:c:c' (x)}$
 $-0.95 * \text{'C-C' (x)}$
 $-0.95 * \text{'C-C-N-C' (x)}$
 $-0.94 * \text{'C-O-C=O' (x)}$
 $-0.94 * \text{'c:c:c:c:c:c-S' (x)}$
 $-0.94 * \text{'c:c:c:c:c-S' (x)}$
 $-0.94 * \text{'c:c:c:c-S' (x)}$
 $- ...)$

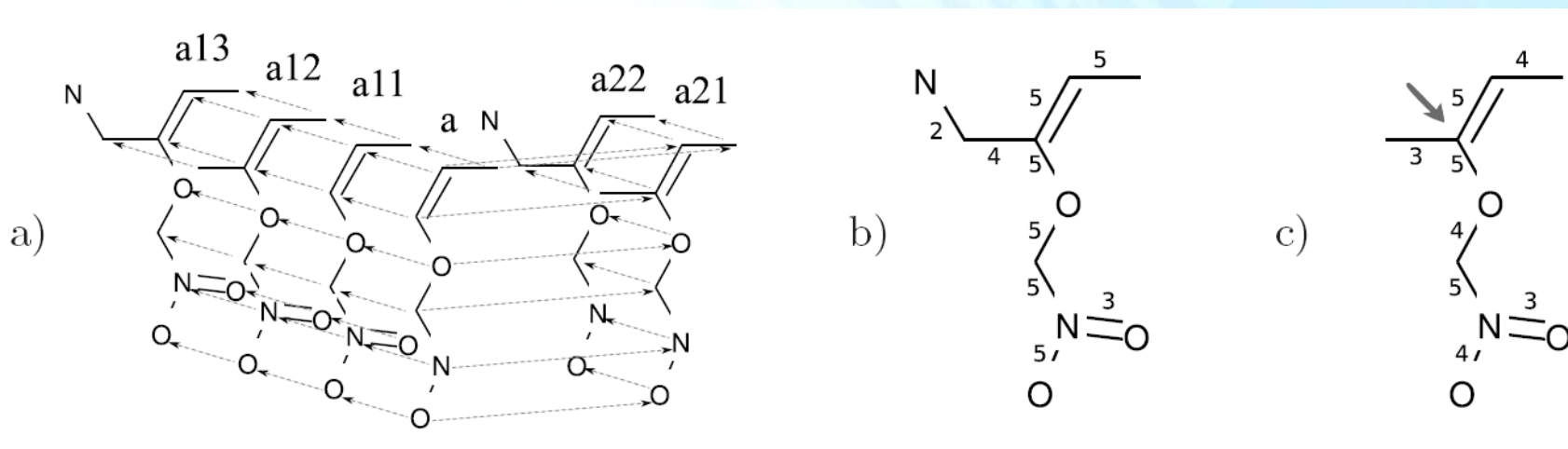
Scalability: A New, Practical Class of Substructures

- A new, practical class of substructures: *backbone refinement classes (BBRC)*, i.e., *trees sharing a common backbone*
- Then pick the most significant representative from this class
- > 23,000 compounds from NCI Yeast Anticancer Drug Screen data: BBRC representatives computed in 4m52s, other approaches did not even finish
- 87,264 possible; producing reasonable coverage of structures



Latent Structure Pattern Mining

- Automatically discovering structural alerts
- 3 steps: (a) align, (b) stack, (c) compress
- Results for: blood-brain barrier, bioavailability, ...



Structural Clustering and Local Models

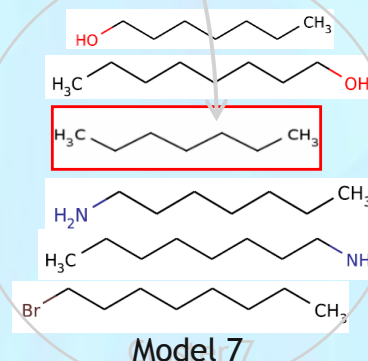
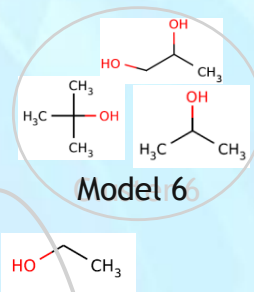
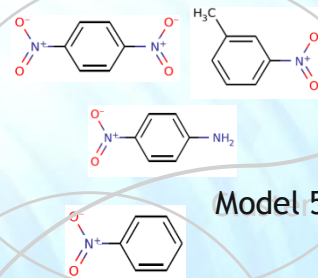
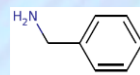
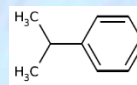
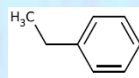
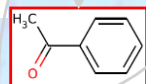
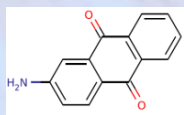
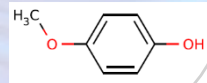
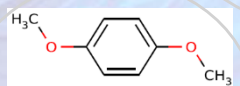
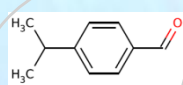
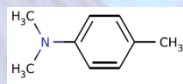
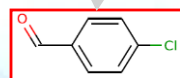
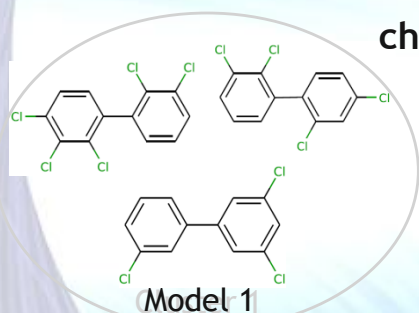
Training set

Test set

preprocessing

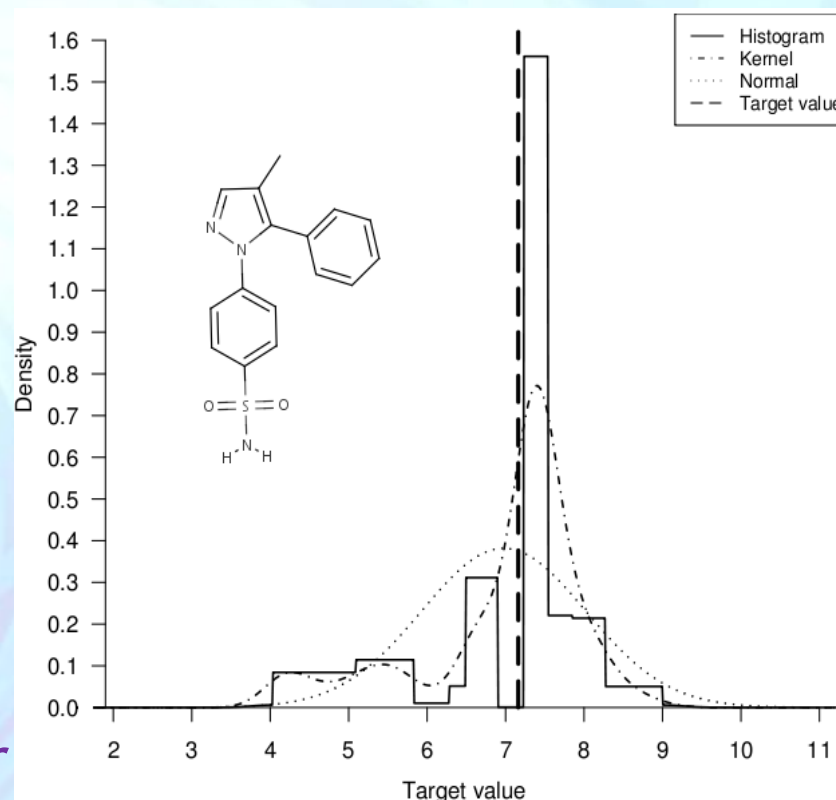
assign to clusters

chemical space



Fast Conditional Density Estimation for QSARs

- Prediction of distribution of activities
 - not point estimates
 - quantifying uncertainty
- Doing it fast...
- ...using general purpose machine learning as plug-in
- Then use histogram estimator, Normal estimator, Kernel estimator



Summary

- Algorithms: descriptor calculation, feature selection, classification (SAR) and regression (QSAR), ...
- Simple API for algorithms
- Development of useful novel algorithms:
 - substructure generation, structural clustering, local models, fast conditional density estimation, multi-label classification, ...

Development and Use of Predictive Toxicology Applications

An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Validation

presented by Haralambos Sarimveis
(National Technical University of Athens, Greece)

Use of QSARs under REACH (Annex IX)

- Acceptance of QSAR results - BOTH positive and negative results will be accepted if
 - Models have been **validated**
 - Models are adequately **documented** and meet acceptance criteria for a given application- “fit for purpose” concept

Compelling Needs of Users

Integrated Testing

in silico

in vitro

TTC

Read
Across

Category
Formation

REACH Reporting
(QPRF, QMRF)

Applicability
Domain

Validation

Human
Data

OpenTox Framework - Standards

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/33/37/37849783.pdf
- QSAR Model Reporting Format (QMRF)
qsardb.jrc.it/qmrf/help.html
- QSAR Prediction Reporting Format (QPRF)
ecb.jrc.it/qsar/qsar-tools/qrf/QPRF_version_1.1.pdf

Reports

REACH

- Guidance on Information Requirements and Chemical Safety Assessment

Part F

- Chemicals Safety Report
- Appendix Part F
guidance.echa.europa.eu/guidance_en.htm

OECD - The Organisation for Economic Cooperation and Development

- Intergovernmental Organisation grouping 30 industrialised countries, aiming to: support sustainable economic growth, boost employment, raise living standards, maintain financial stability, assist other countries' economic development, contribute to growth in world trade.
- The OECD works on global issues in different areas, such as economy, society, governance, development, finance, innovation, sustainability.
- In November 2004, the OECD member countries agreed on the **principles for validating (Q)SAR models** for their use in regulatory assessment of chemical safety. The agreed principles provide member countries with basis for evaluating regulatory applicability of (Q)SAR models and will contribute to their enhanced use for more efficient assessment of chemical safety.

	OECD Principle	To facilitate the consideration of a (Q)SAR model for regulatory purposes, it should be associated with the following information::
1	Defined Endpoint	Ensure clarity in the endpoint being predicted by a given model, since a given endpoint could be determined by different experimental protocols and under different experimental conditions.
2	Unambiguous Algorithm	Ensure transparency in the model algorithm that generates predictions of an endpoint from information on chemical structure and/or physicochemical properties. The issue of reproducibility of the predictions is covered by this Principle.
3	Defined Applicability Domain	(Q)SARs are inevitably associated with limitations in terms of the types of chemical structures, physicochemical properties and mechanisms of action for which the models can generate reliable predictions
4	Goodness-of-fit, robustness and predictivity	Internal performance of a model (as represented by goodness-of-fit and robustness) and the predictivity of a model (as determined by external validation).
5	Mechanistic interpretation (if possible)	The intent of Principle 5 is not to reject models that have no apparent mechanistic basis, but to ensure that some consideration is given to the possibility of a mechanistic association between the descriptors used in a model and the endpoint being predicted.

	OECD Principle	OpenTox addresses Validation Principles by...
1	Defined Endpoint	providing a unified source of well defined and documented toxicity data with a common vocabulary
2	Unambiguous Algorithm	providing transparent access to well documented models and algorithms as well as to the source code
3	Defined Applicability Domain	integrating tools for the determination of applicability domains during the validation of prediction models
4	Goodness-of-fit, robustness and predictivity	providing scientifically sound validation routines for the determination of errors and confidences
5	Mechanistic interpretation (if possible)	integrating tools for the prediction of toxicological mechanisms and the recording of opinions and analysis in reports

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).
- 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>

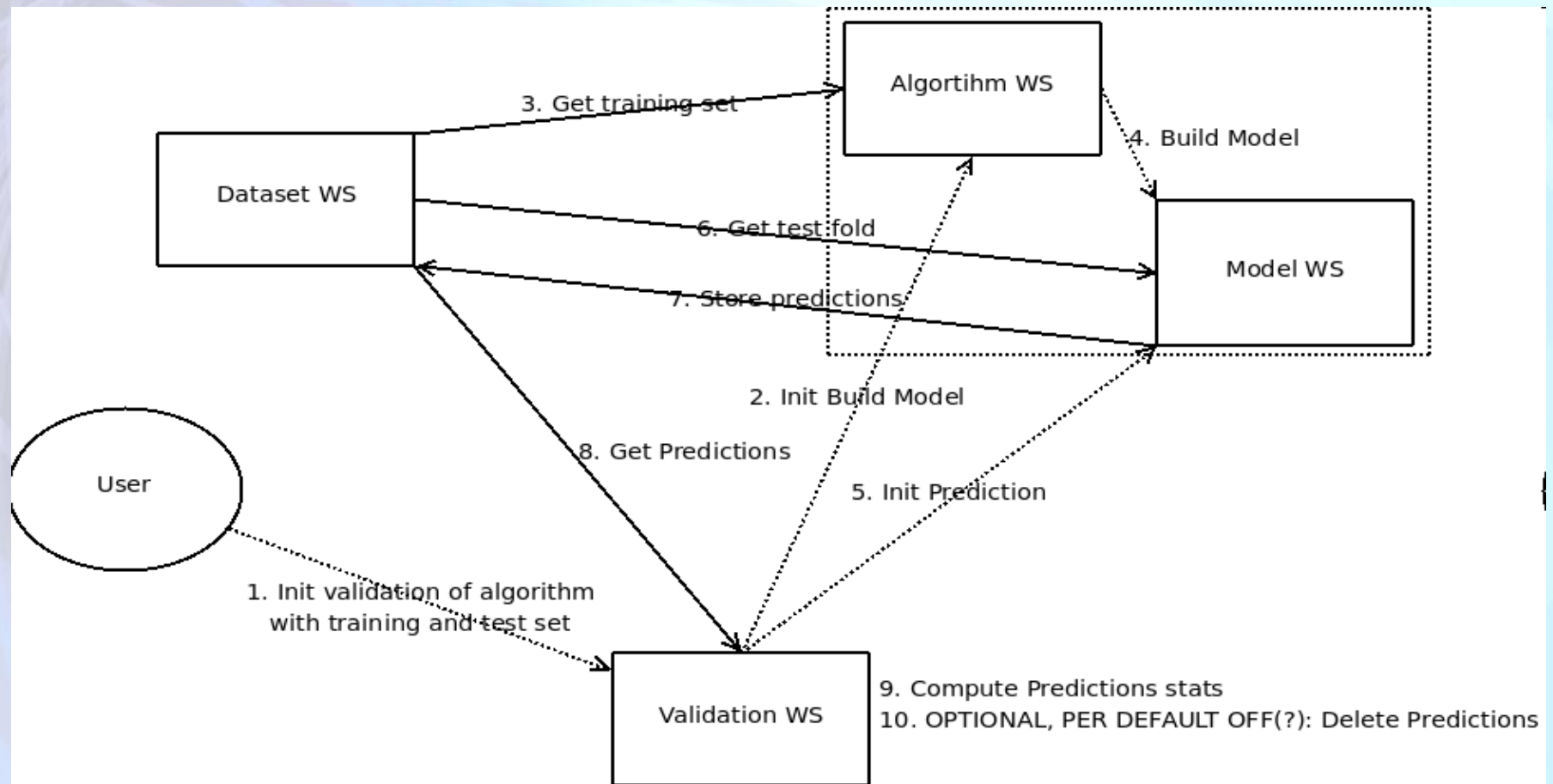
Validation API

Description	Method	URI	Parameters	Result	Status codes
Get all validations	GET	/	-	List of validation URIs	200,404
Retrieves a validation representation	GET	/ {id}	-	Validation representation in one of the supported MIME types	200,404
Validates a model on a test dataset	POST	/	model_uri test_dataset_uri test_target_dataset_uri (default = test_dataset_uri)	Validation URI or Task URI	200,400,404,500
Builds a model on a training dataset and validates it on a test dataset	POST	/	algorithm_uri prediction_feature algorithm_params (string, default="") training_dataset_uri test_dataset_uri test_target_dataset_uri (default = test_dataset_uri) y_scramble (boolean, default=false) y_scramble_seed (integer, default=1)	Validation URI or Task URI	200,400,404,500
Splits a dataset into training and test dataset according to a certain ratio, and performs a validation	POST	/training_test_split	algorithm_uri prediction_feature algorithm_params (string, default="") dataset_uri split_ratio(float, default=0.66) random_seed(integer, default=1) y_scramble (boolean, default=false) y_scramble_seed (integer, default=1)	Validation URI or Task URI	200,400,404,500
<i>OPTIONAL:</i> Performs a bootstrap validation	POST	/bootstrap	algorithm_uri prediction_feature dataset_params (string, default="") dataset_uri bootstrap_percentage(float, default=0.66) random_seed(integer, default=1) y_scramble (boolean, default=false) y_scramble_seed (integer, default=1)	Validation URI or Task URI	200,400,404,500
Deletes a validation.	DELETE	/ {id}	-	-	200,404

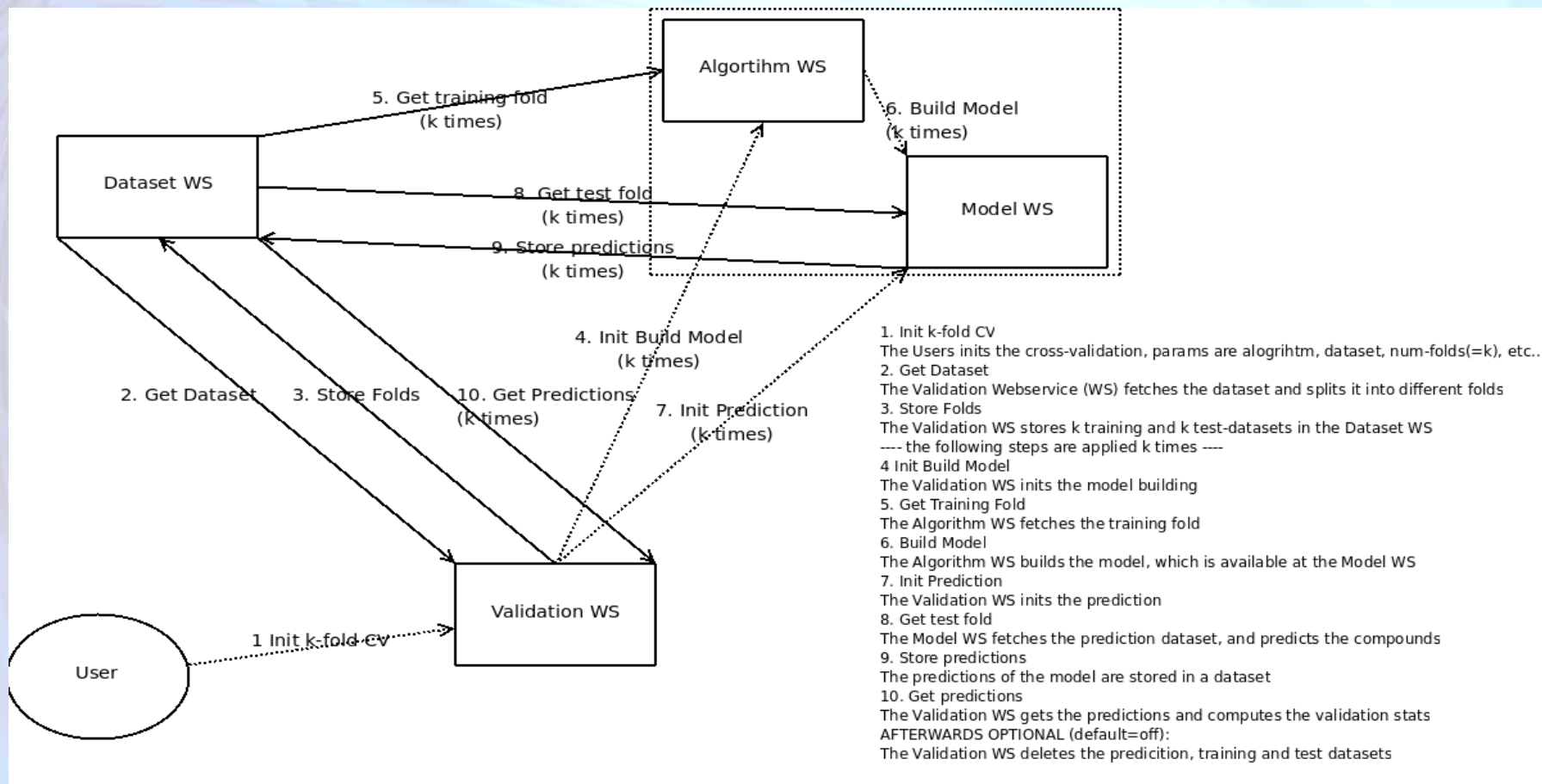
Cross-Validation API

Description	Method	URI	Parameters	Result	Status codes
Get all cross-validations	GET	/crossvalidation	-	List of crossvalidation URIs	200,404
Retrieves a cross-validation representation	GET	/crossvalidation/{id}	-	Cross-Validation in one of the supported MIME types	200,404
Returns all (k) validations that belong to a crossvalidation	GET	/crossvalidation/{id}/validations	-	List of validation URIs	200,404
Performs a k-fold cross-validation.	POST	/crossvalidation	algorithm_uri prediction_feature algorithm_params (string, default="") num_folds (integer, default=10) random_seed (integer, default=1) stratified (boolean, default=true) y_scramble (boolean, default=false) y_scramble_seed (integer, default=1)	Cross-Validation URI or Task URI	200,400,404,500
Performs a leave-one-out cross-validation.	POST	/crossvalidation/loo	algorithm_uri prediction_feature algorithm_params (string, default="") y_scramble (boolean, default=false) y_scramble_seed (integer, default=1)	Cross-Validation URI or Task URI	200,400,404,500
Deletes a cross-validation.	DELETE	/crossvalidation/{id}	-	-	200,404

Validation WorkFlow



Cross-Validation WorkFlow



Applicability domain/confidence in prediction

- A definition of AD which is also used by the OECD is the following: “The applicability domain of a (Q)SAR model is the response and **chemical structure space** in which the model makes predictions with a given **reliability**.”
- Furthermore, OECD advises that the AD principle should be applied in a **model-specific manner**. Thus, every model should be associated with its own AD derived not only on the chemicals in the training set but also on the descriptors and (statistical) approach used to develop the model. Ideally, the AD should be defined and documented by the model developer.
- Related to the concept of an AD is the concept of **confidence** in predictions inherent in some learning algorithm, so that the predictive model itself provides estimation of applicability domain. For example, classification algorithms do not only provide a categorical class label, but also a **probability** with which the class is predicted. The main difference is that the confidence is only known when the model is already applied, whereas the applicability domain is defined on the input space directly.

Implemented 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

Future Work

- Validation with a test set that is **completely unknown** to the model developer is certainly the gold standard in this area, because there is no way to cheat voluntarily or involuntarily (e.g. by "optimizing" model parameters for a specific test set).
- OpenTox will provide facilities to access **confidential** (inhouse) data. For validation purposes we will provide a facility to test (Q)SAR models remotely against confidential datasets without getting access to the actual entries of the database to ensure security and confidentiality of proprietary data.
- Confidential validation data will be sought from **external sources** including members of the advisory board.
- OpenTox already provides facilities to protect confidential information located at **URIs**. Two tasks are involved here:
 - **Authentication**: Confirming the identity of the user requesting access
 - **Authorisation**: Granting the confirmed identity access according to a set of restrictions described in policies

Development and Use of Predictive Toxicology Applications

An OpenTox Workshop
19 Sep 2010, Rhodes, Greece

Reporting

presented by Andreas Karwath
(Albert-Ludwigs-Universität, Freiburg, Germany)

Compelling Needs of Users

Integrated Testing

in silico

in vitro

TTC

Read
Across

Category
Formation

REACH Reporting
(QPRF, QMRF)

Applicability
Domain

Validation

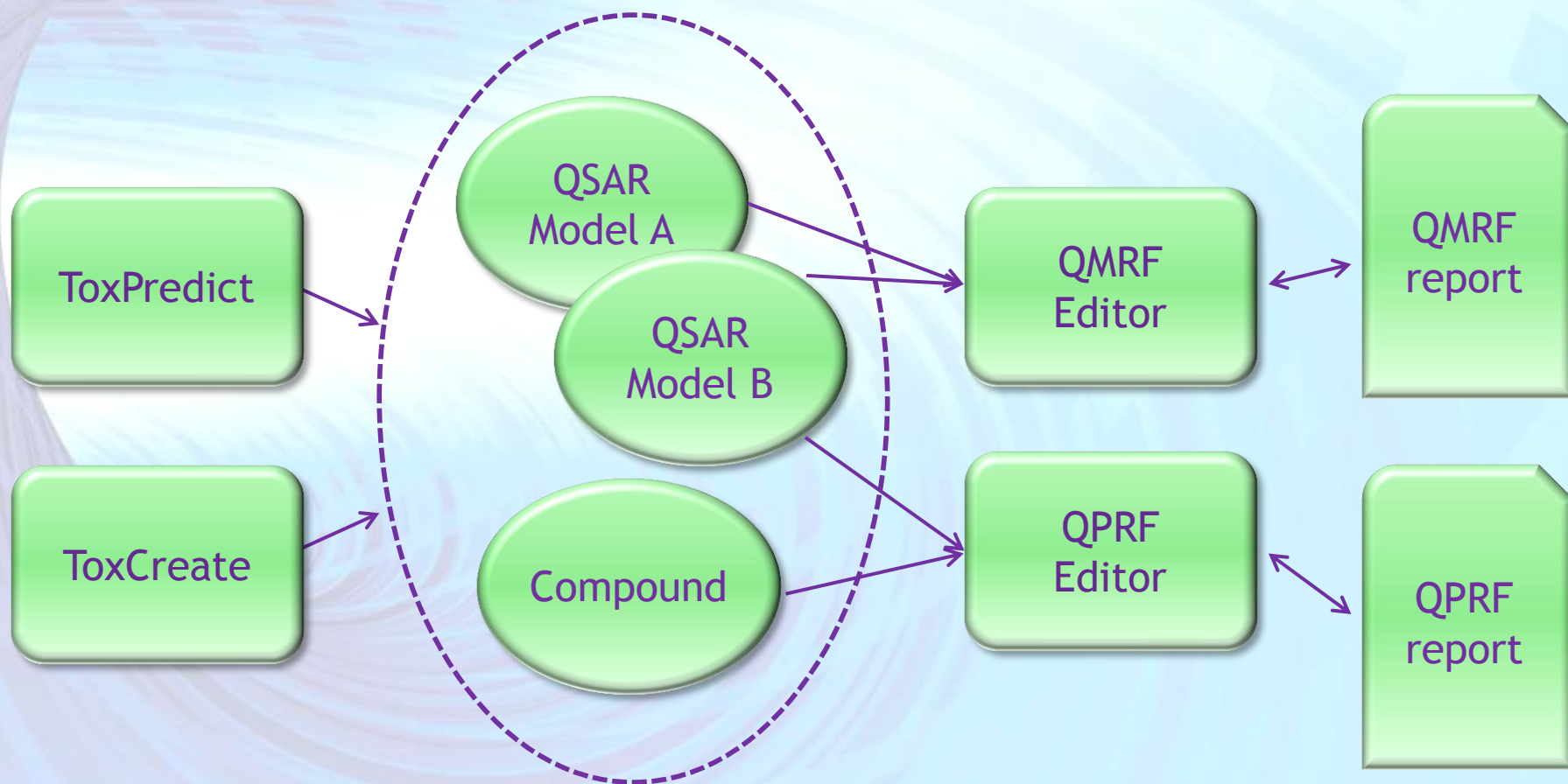
Human
Data

REACH reporting formats

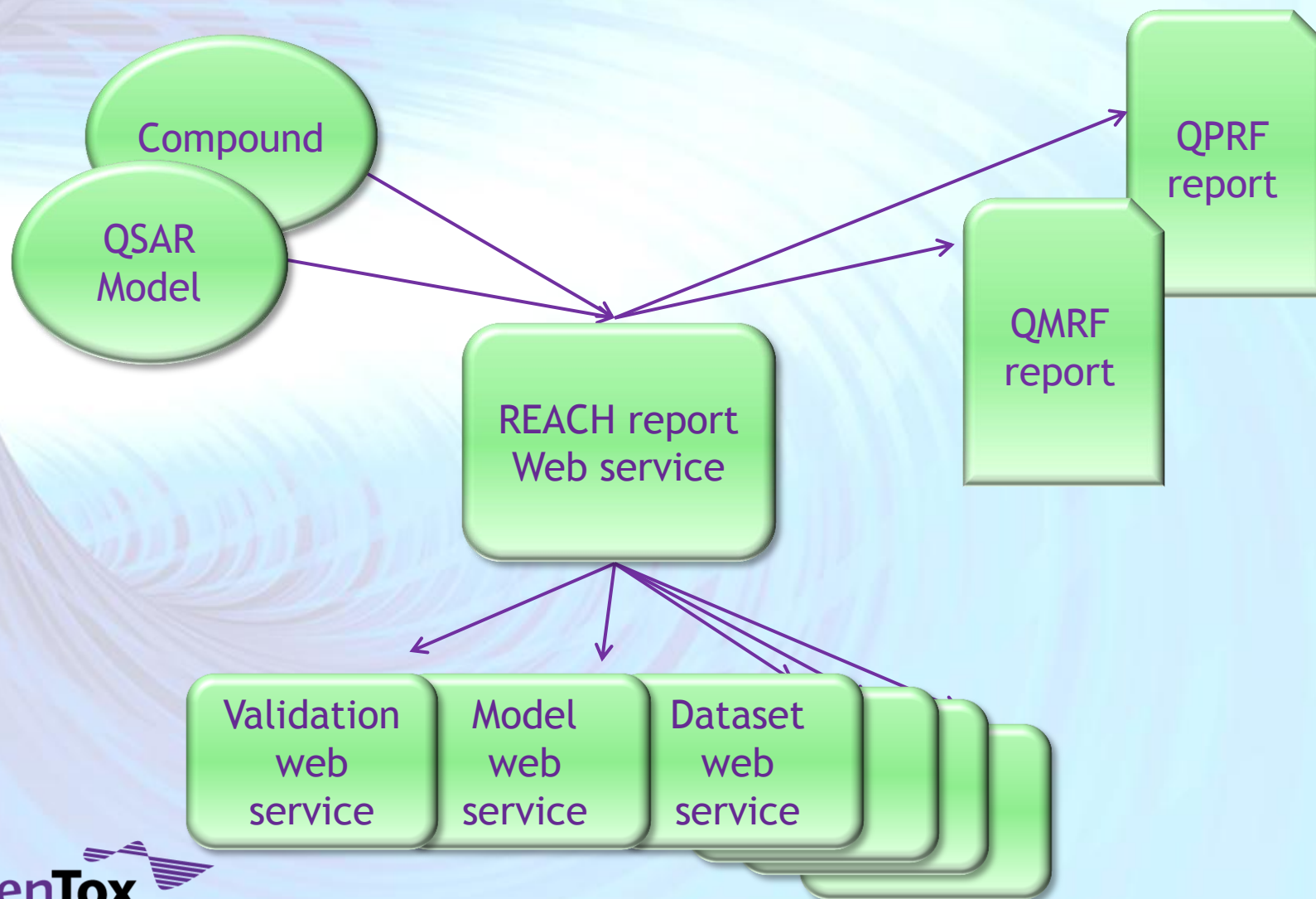
QMRF and QPRF :

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

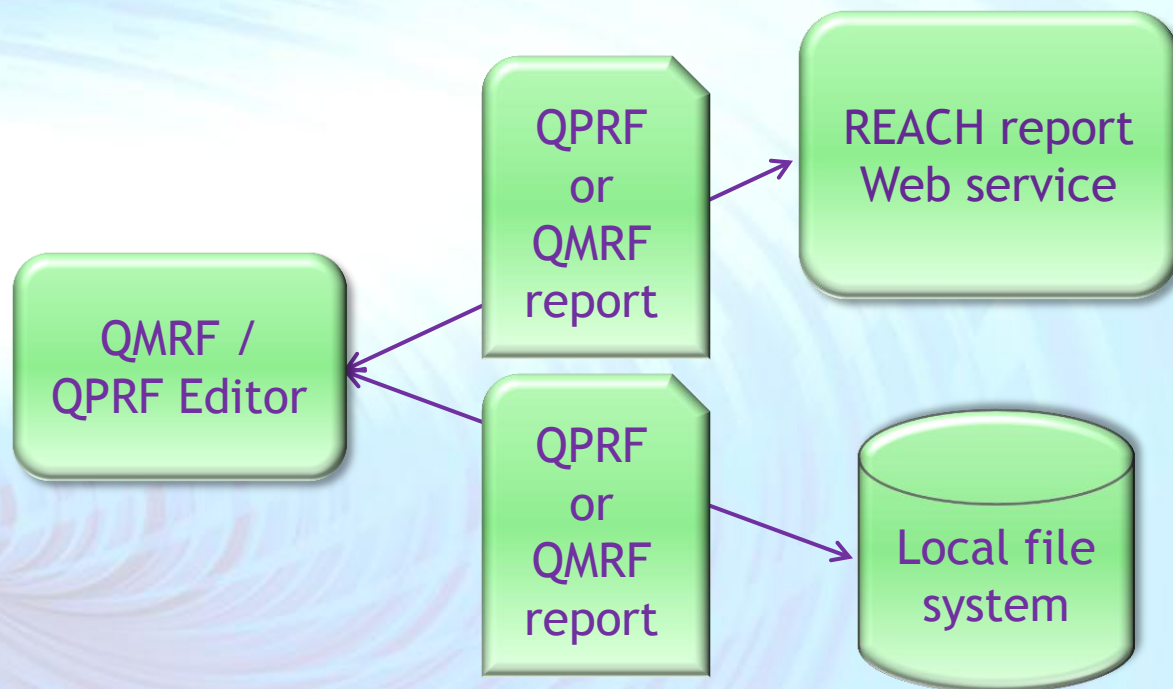
User Perspective



Creating Reports



Storing and Editing Reports



Implemented Reporting Services - QMRF

Report web service:

- Automatic generation of reports, including:
 - Meta-information (creation date, algorithm, model endpoint, ...)
 - Model training data
 - Validation results (cross-validation, bootstrapping, ...)
 - Prediction results on external test-data
- Reports can be downloaded/uploaded/deleted

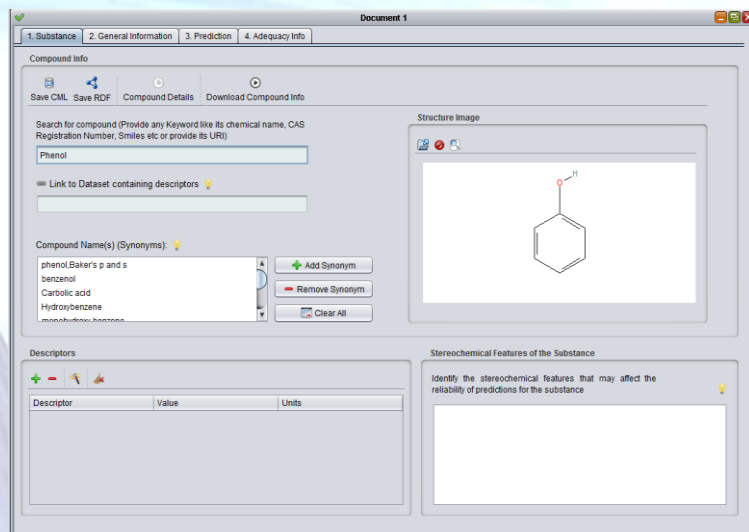
Implemented Reporting Services - QMRF [2]

QMRF Editor:

- Based on existing and EU approved implementation (see <http://qsardb.jrc.it/qmrf>)
- Comprehensive functionalities (edit/store reports, export to pdf)
- Extended to communicate with web service (download and upload reports)
- Embedded into ToxCreate

Implemented Reporting Services - QPRF

QPRF Editor (Q-Edit):



...Using AMBIT services

- Version 0.1.3 (alpha), heading towards the first beta version
- PDF Creation fully supported
- Compound lookup services facilitate users to find the compound they are looking for
- Similarity Search using the Q-Edit GUI
- Reports are (locally) stored in a binary format (RDF is under development)
- Available for download from <http://github.com/alphaville/Q-edit>

Future Work

QMRF Report Services and Editor:

- Include more automatically generated information:
 - Detailed description of model and algorithm
 - Related models
 - Authors
 - ...
- Enable Authentication

QPRF Report Services and Editor:

- Establish web services for QPRF
- Enable Authentication

Final words...

For more information, visit

www.opentox.org

Contact Project Coordinator:

barry.hardy@douglasconnect.com

+41 61 851 0170

**We welcome your
involvement!**



OpenTox - An Open Source Predictive Toxicology Framework, www.opentox.org, is funded under the EU Seventh Framework Program: HEALTH-2007-1.3-3 Promotion, development, validation, acceptance and implementation of QSARs (Quantitative Structure-Activity Relationships) for toxicology, Project Reference Number Health-F5-2008-200787 (2008-2011).