# Introduction to SAR/QSAR analysis: On-line chemical database and modelling environment (OCHEM)

Dr. I.V. Tetko, BIGCHEM GmbH and Institute of Structural Biology, Helmholtz Zentrum München - German Research Center for Environmental Health (GmbH)

Public platform: www.ochem.eu
Online manual:  http://docs.ochem.eu/display/MAN

HelmholtzZentrum münchen
German Research Center for Environmental Health

# OCHEM tutorial handout

**Table of Contents**

# 1.  General concepts

In this chapter, we will learn the general concepts of the OCHEM interface
How (and why) to register a new account
What are the basic design components of OCHEM
What are the basic elements of OCHEM
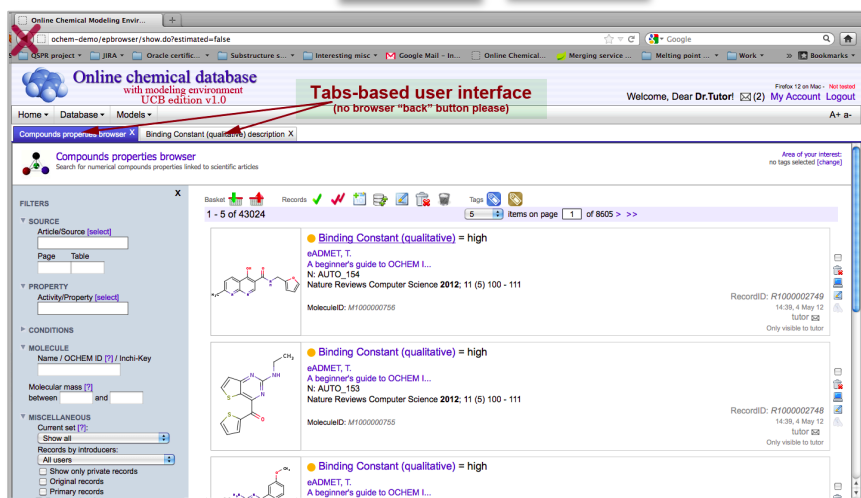
# OCHEM tutorial handout

## 1.1    Before we start



The OCHEM installation is running on servers in HMGU. It can be installed on a virtual machine or on computers for better performance.

**System requirements**: For optimal performance, the host machine should have at least 8 CPU cores, 16GB RAM and about 100GB disc space.

# OCHEM tutorial handout



OCHEM is a web-based platform. Users access it with a simple web browser, similar to the way they access services like Gmail or Facebook.

## How to access OCHEM?

The public version of OCHEM is available at www.ochem.eu, but an in-house installation can be run inside a company/University and be accessible inside of the intranet only.

## Which browser to use?

To get the best experience, it is recommended to use the latest Firefox browser. Chrome and Safari are also supported. Unfortunately, at the moment OCHEM **does not fully support Internet Explorer (improves with release of new versions)** and Konqueror Web browser.

## User interface remarks.

OCHEM uses "tabs" extensively in the user interface. New dialogues are often conveniently opened in new tabs, as it is shown on a screenshot on the left.

Please note, that the browsers "back" button is not compatible with the tabbed interface.

## 1.2    User account creation and user login

### Please, login

**Instant login**

In order to access OCHEM, you must login. If you do not wish to register now, you can login as a guest. You can also use your Facebook account to login.

LOGIN AS A GUEST    LOGIN WITH facebook

**Already have an account?**

If you already have an account, please enter you login and password below:

Login ID

Password

LOGIN    PASSWORD REMINDER

**Join OCHEM - register a new user!**

Create a free account to upload data, create and apply QSAR models, screen chemical libraries and many more. Registered users can correct data uploaded by other registered users publish models. As a registered user, you can configure flexible access policies for your data and models.

REGISTER A NEW USER

In order to use the OCHEM web platform, a user has to register an account and login using this username and password.

Users can login as a guest user (with limited privileges) or as a registered user.

**1.** The first option on the login interface is the instant login (without registration). You can log in as a guest user. For this tutorial, you are encouraged to register an account.

**2.** If you have already created an OCHEM account, you can login using your username and password.

**3.** The third option is to register a new account. Please, register an account if you have not done so yet.

**User account**
Details of your personal OCHEM account

**Registration Information**

Login* ▢ *(min. 4 characters and max. 20 characters)*
CHECK AVAILABILITY

e-mail* ▢

Password* ▢

Confirm password* ▢

**Personal Information**

Title* -- please select --  *Please, select a form of address!*

First name* ▢

Last name* ▢

Affiliation ▢

Form of organization* -- please select --  *Please, select a form of organization!*

City ▢

State ▢

Country ▢

Zip ▢

Phone ▢

Occupation ▢

Company ▢

WebSite ▢

In order to create a new OCHEM user account, the following information is required:

**4.** Choose your login name and check if it is available (i.e. it is not yet used by another user). Login names should be at least 4 characters long and not longer than 20 characters.

**5.** A valid e-mail address is required for the automatic notification system.

**6.** Furthermore a password for the acount should be chosen and confirmed.

**7.** Additional personal information like academic title, first and last name of the user and the form of organization this person is working in is required to finish the user account creation

Note:
   Registered users have access to more features than guest users (i.e., can upload data and develop models).

   If users provide detailed information about themselves, their account will **be validated** by the OCHEM administrator. This will allow the users to run larger tasks, export more data and edit data of other validated users.

## 1.3    Data browsers



An important user interface element in OCHEM is a browser. OCHEM has various browsers for all kinds of database entities.

**Main browsers include:**

- Experimental property browser (record browser)
- Molecules browser
- Properties browser
- Conditions browser
- Units browser
- Articles browser
- Journals browser
- Baskets browser
- Tags browser
- Models browser

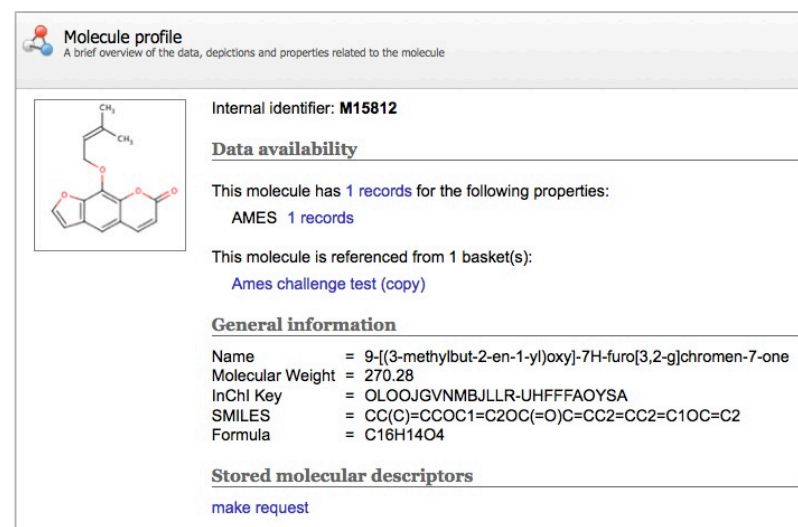**Main elements of a browser:**

1. Items (e.g., data records, models, articles, tasks)
2. Page bars – can be used to navigate items
3. Filters – can be specified to narrow down the displayed items to a specific area
4. Global toolbars – can be used either to manipulate (delete, modify) several items simultaneously, or create new items
5. Item toolbars – can be used to perform operations on a specific item

## 1.4 Item profiles



**Exemplary item profiles:**
- Compound property editor (record editor)
- Molecule profile
- Article profile
- Dataset ("basket") profile
- Property editor
- Unit editor

# 2. Using OCHEM

In this chapter, we will how to use OCHEM for several typical scenarios, including

- Search of properties using compound properties browser
- Grouping and exporting records
- Use of ToxAlert  for data exploration
- SetCompare Tool for basket comparison
- Prediction of properties for molecules using models

## 2.1    Compound properties browser



Compound properties browser is one the main dialogues in OCHEM. It allows you to browse experimental data records using a variety of filters.

Please, try using different filters, e.g.:

1. Show only the data for the "Ames" property
2. Filter the records by substructure
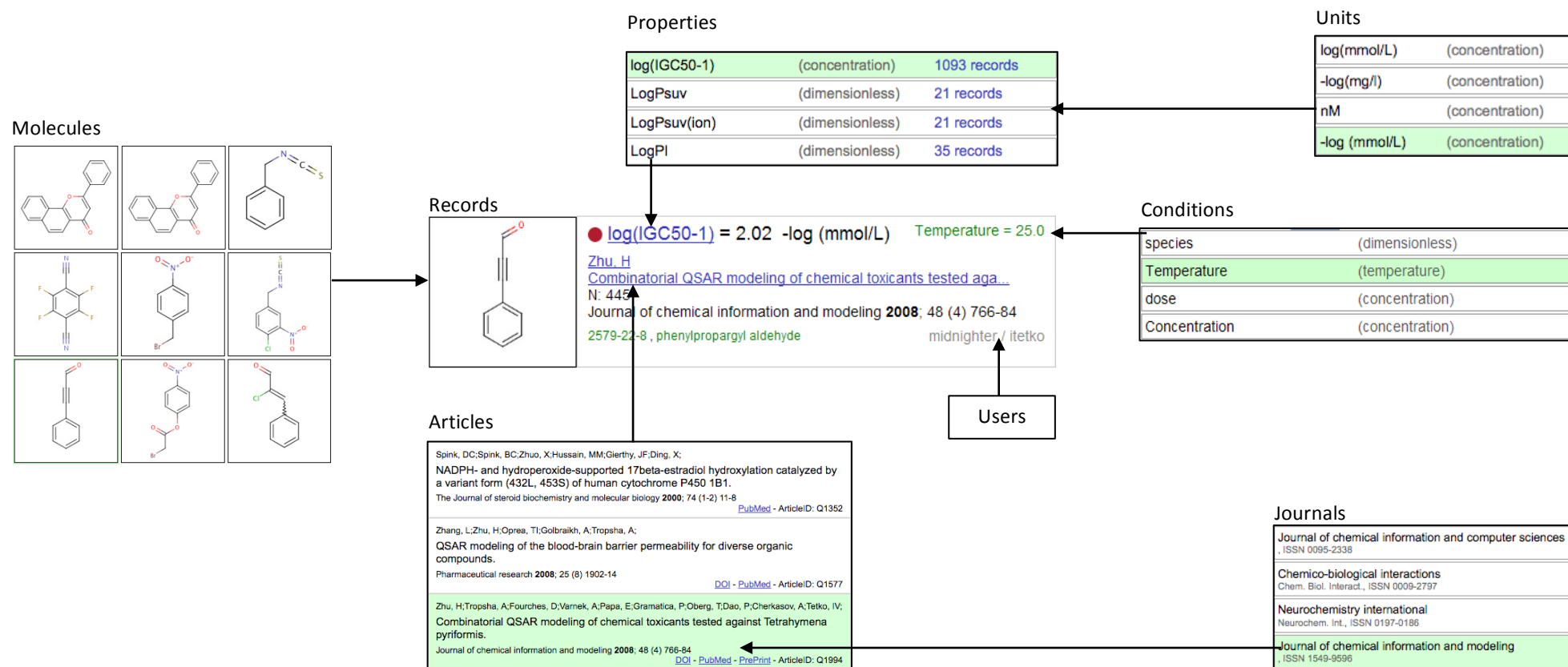3. Sort the data by ascending molecular weight

For further use, records can be selected individually ☑ or, more commonly, it is possible to select all the records matching current filters ✔

The selected records can be put into a dataset (so called "basket"), which later can be reused for any kind of tasks, e.g. for development or validation of QSAR models.

The selected records are shown at the top of the left panel and persist until they are deselected (cleared). Clear them before selecting new records unless you would like to add them to the previous records.

**Other filters of property records are available from browsers of Articles, Properties, Baskets, Models, etc.**

## 2.2 Data structure



Experimental property (or "record") – a value for a property for a specific molecule published in a specific article or book.

This means that:

One molecule can have multiple records associated with it (measurements for different properties, measurements for the same property published in different articles, etc.)

One article can hold multiple records for multiple properties for multiple molecules

Most of essential OCHEM operations (such as QSAR modeling) are performed on datasets of records (and not molecules)

## 2.3    Working with baskets and data export



### 1. Finding interesting records

Filtering of records using filters such as article, model, property, basket or selection of records one by one, on page or all filtered
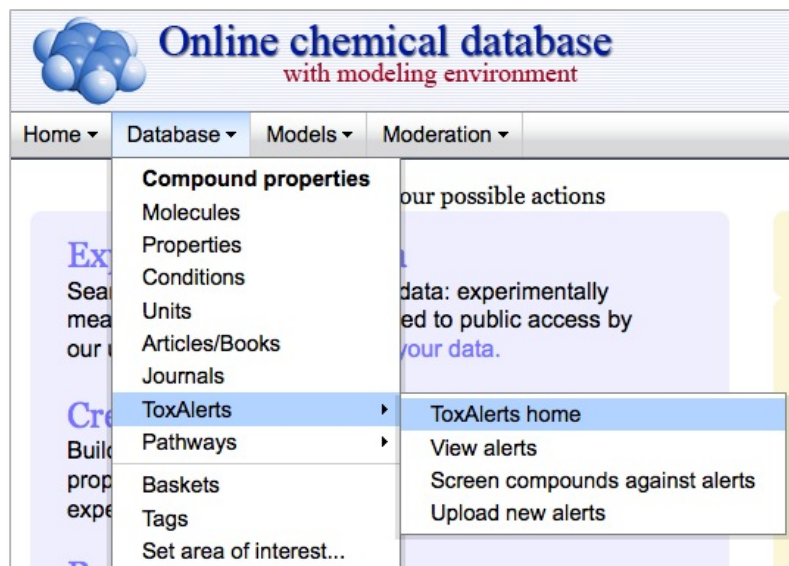


### 3. Exporting the records



### 2. Exploring the baskets

## 2.3 ToxAlert utility



Alerts are structural features that are known to be associated with a particular activity. Typical structural alerts might indicate carcinogenicity or general toxicity. Screening molecules against alerts can be simple and, most importantly, easily interpretable

The *ToxAlert* utility allows one to screen a set of molecules against a set of structural alerts. OCHEM comes with thousands of alerts for a number of endpoints. It is also possible to introduce your own alerts

1. To get to the ToxAlert utility just select it from the menu bar.



2. A welcome page is the entry point for further actions, like
   - overview available alerts
   - upload new structural alerts
   - screen molecules against structural alerts

**3.** The structural alerts browser gives an overview of the available alerts in the system.

- Existing alerts can be filtered by their category
- New alerts can be uploaded
- Selected alerts can be used in a screening against a set of structures



**Screening compounds against alerts**

**4.** Select the compounds you would like to screen. This can be a prepared basket from the OCHEM platform, a single structure drawn or automatically fetched by its name, or like in this example an uploaded file (SDF, Smiles, Excel).

**5.** Optionally, you can screen against alerts for a particular endpoint or alerts from a particular publication.

**6.** As result you can see the structures for which a particular alert was found and from which publication this alert stems.

Further information about ToxAlert screening is available in the OCHEM documentation (ToxAlerts)

## 2.5    Set Compare utility



**The SetCompare** utility allows juxtaposing two sets of chemical compounds and finding distinguishing features of each set. For example, you can compare active and inactive compounds for a particular property.

The utility helps to address the following questions:
- What are the distinguishing structural features of active compounds?
- How significant are these results statistically?
- Which are the compounds that possess a particular important structural feature?

SetCompare is accessible from the menu bar.

In the first page of the wizard two sets have to be selected. With the set comparison utility two sets can be examined with respect to common structural alerts and common descriptors.

## 2.6 Model application



**Getting to model application**

**1.** Open the model applier browser from Models > Apply a model.

**Models applier browser**

- The models applier browser lists all the models (public and private, developed by the user). It shows the model name, the predicted properties, used training set, used machine learning method and the creation date.
- The [icon] icon links to model export.



**2.** Please, find the model you want to apply and click "apply the model".

Note:
You can check and apply multiple models simultaneously.

**Application of regression model:**

**3.** Now it is time to provide the compounds you would like to predict. There are several possibilities:
   - Upload structures (e.g., in SD-format)
   - Provide SMILES for a single molecule
   - Draw a structure in a visual structure editor
   - Use an earlier created basket
   - Select a certain set of records / molecules by a tag

Selecting the prediction scenario and disabling the prediction cache are additional options for the prediction process.

**4.** Click on next button to start the application
Wait until the calculations are completed. Again, if the task is taking long time, it is possible to fetch results anytime later from the pending tasks browser (Models > Pending tasks menu).

**5.** Now the calculations have finished. The results include prediction values, and accuracy estimates for each predicted compound. Additionally, there is an estimate for the overall prediction accuracy for the set.

**6.** Predictions can be exported to an Excel or CSV sheet, SDF files or R scripts.

**7.** Prediction results of each single structure are listed in the browser. There are the predicted values themselves, distance to model value and an estimation of the accuracy (RMSE), as well as the originally measured value if it is known.

**Application of classification models**

For application of classification models, the same steps have to be done as for a regression model.

**8.** The result browser shows the predicted class together with an estimation of the accuracy. In this example case, the model was applied to a single drawn structure.

# 3 Modeling framework

In this chapter, we learn how to upload data and develop QSAR models

Development of multiple models using Comprehensive Modeling (CM)
Development of single models

The basic steps of a QSAR modeling lifecycle: prepare data, configure model, train the model, analyse results and use the model to predict new compounds

# OCHEM tutorial handout

## 3.1 Comprehensive modeling



1. The "comprehensive modelling" feature accessible via the "Models" menu is an advanced feature that allows you to easily create multiple models based on different descriptor sets and training methods.

   With this feature, you can create dozens of models simultaneously and directly compare their performance.



2. In the following dialog, first you should select your training set.

3. You can see a set of predefined configuration templates for several training methods, molecular descriptors, descriptor selection methods and model validation.

   The checked methods will be applied using the "all against all" principle. On the following screenshot, we selected methods, descriptor sets, descriptor selection method and validation method, which results into six models.

   We selected only six models for the reason of speed. Normally, you can run dozens or hundreds of models, depending on available calculation resources.

4. Now we are ready to launch all models.

**5.** Please, wait until OCHEM starts the necessary calculation tasks.



**6.** When done, you are forwarded to the success page, from which you can directly go to the **models summary page**.



**7.** The models summary page built for a particular basket is also available via Basket browser (menu Database > Baskets), by clicking ⊞ icon for your basket.

**8.** The models summary page shows all the models (ready and pending) for the selected basket.

The models are grouped by methods, descriptors and validation protocols. Currently, we see that our four models are still running.

You can return to this page at any time to check the status of your models or click "refresh" to update the dialog. Normally, the creation of multiple models takes a while.



**9.** To calculate statistics for all the completed models, press the "fetch statistics for ready models".

Multiple models overview

Predicted property: log(IGC50-1)
Training set: T. pyriformis train

Metrics [ RMSE - Root Mean Square Error ] for [ Training set ] Validation: [ Cross-Validation (18 models) ]

| | ASNN | ASNN(2) | ASNN(3) | LibSVM | FSMLR | PLS |
|---|---|---|---|---|---|---|
| OEstate | 0.44 | + | + | + | + | + |
| StructuralAlerts | + | 0.76 | + | + | + | + |
| EState | + | + | 0.45 (+1 models) | + | + | + |
| CDK (constitutional, topological, geometrical, electronic, hybrid) | + | + | 0.49 | 0.64 | 0.58 | 0.58 |
| ALogPS, OEstate | + | + | 0.54 | 0.56 | 0.65 | 0.56 |
| ChemaxonDescriptors (7.4) | + | + | 0.5 | 0.69 | 0.75 | + |
| EState, ALogPS | + | + | 0.48 | + | + | + |
| GSFrag | + | + | 0.62 | + | + | 0.73 |

Refresh

Export as Excel file
Export as R script

**10.** We can see all four of our models are ready. The numbers in the cells ("metrics") show the root mean square error.

In this particular case, we can immediately observe that neural network models (ASNN) have lower errors (RMSE) than e.g. the PLS models.



Multiple models overview

Predicted property: log(IGC50-1)
Training set: T. pyriformis train

Metrics [ RMSE - Root Mean Square Error ] for [ Training set ] Validation: [ Cross-Validation (18 models) ]

RMSE - Root Mean Square Error
MAE - Mean Absolute Error
R2
Q2
Model size

| | ASNN | ASNN(2) | ASNN(3) | LibSVM | FSMLR | PLS |
|---|---|---|---|---|---|---|
| | 0.44 | + | + | + | + | + |
| | + | 0.76 | + | + | + | + |
| EState | + | + | 0.45 (+1 models) | + | + | + |
| CDK (constitutional, topological, geometrical, electronic, hybrid) | + | + | 0.49 | 0.64 | 0.58 | 0.58 |
| ALogPS, OEstate | + | + | 0.54 | 0.56 | 0.65 | 0.56 |
| ChemaxonDescriptors (7.4) | + | + | 0.5 | 0.69 | 0.75 | + |
| EState, ALogPS | + | + | 0.48 | + | + | + |
| GSFrag | + | + | 0.62 | + | + | 0.73 |

Refresh

Export as Excel file
Export as R script

**11.** It is also possible to display other statistical parameters, such as R2 or Q2, using the drop-down box.

**12.** You can perform row-wise or column-wise batch operations, e.g., delete the models or create new models.

**13.** You also can create new models individually by pushing "+" sign in the "missing" cells.

## 3.2 Development of a single model

### Select training set, machine learning method and internal validation options





To start the model development process, open "Model > Create a model" from the menu panel. The first page of the model creation "wizard" asks you to select a training set, external validation sets (optional), a machine learning algorithm and an internal model validation technique.

1. Select the training set and optionally one or more external validation sets that you have prepared before by clicking on the [...] label and the "Add a validation set" link.

2. OCHEM supports two dozen state of the art machine-learning methods. For this tutorial, we will use defaults for most of the configurable options. Thus, we will select associative neural networks (ASNN) to train the model.

3. You can choose between n-fold cross-validation, bagging and no validation at all. A 5-fold cross-validation is most commonly used.

4. Model configurations can also be imported from earlier model building processes to follow the same protocol.

The model creation process is organized as a "wizard" guiding you through the model configuration process. So click "Next" to navigate forward.

Pre-processing of the molecules includes four options: standardization of some chemical groups for consistency, neutralization of ions, removal of salts and cleaning of given meta-information in the structure file.

5. We will use the default recommended configuration and employ all the available pre-processing options.

## Configure molecular descriptors

Selection of molecular descriptors is an important step that can significantly contribute to the quality of the model.

**6.** For this tutorial, we will use the default selection – E-State descriptors and ALogPS.

Several descriptors and descriptor packages are available on the OCHEM platform, ranging from simple 1D to sophisticated 3D descriptors. If a 3D descriptor is selected, a structure optimization method can be selected in the next step (not shown here).

Furthermore, the output of already existing models can be used as input for the new model to train. I.e. a predicted logP value (if not available) can be used as a molecular descriptor. This functionality is referred to as "feature nets".

The next dialog allows filtering out redundant and correlated descriptors.

**7.** Again, for the purpose of this tutorial, we will use the default values, which include simple filters like pairwise decorrelation.

It is also possible to select the list of desired descriptors manually (advanced).

# OCHEM tutorial handout

## Configure the training method and start calculations



Each machine learning method (e.g., neural networks in our case, KNN, MLR, PLS, etc.) requires additional configuration options. For neural networks, we can configure the training algorithm, the number of neurons, learning iterations and the number of networks in the ensemble.

**8.** We will not experiment here now and will use the default options, which are often a good starting point.



**9.** Finally, we are ready to start calculations. Before starting, please provide the name for your future model.

**10.** Specifying the priority of the calculations is optional and defaults to "normal".

**11.** Please click "start calculations" to start the model training process.

## Distributed model calculation

A waiting-screen that shows you the status of the calculations.
The training process is automatically distributed to several internal calculation units, but still for large datasets it can take a while to complete (from minutes to weeks).

**12.** Although we could have waited, we will opt to click "fetch result later" to get an overview of currently submitted tasks to the system.

The next screen is the list of currently pending tasks, also accessible from menu "Model > View pending tasks". This list displays all tasks that are currently running on the system or have been finished, but not yet fetched by the user.

**13.** Here you can observe the status, terminate running tasks or fetch ready tasks. Please, click "refresh" or check the box for Refresh every minute" to actualize the page.

**14.** When the task has finished, please click the green check button or the model name link to fetch the model and investigate the statistics of the model.

## Save your model



If the calculation was successful, you can see the profile of the ready model. Before saving the model the profile can be investigated further:

**15.** The model profile shows information about the training configuration (used descriptors, machine-learning method and predicted property. It shows the training process statistics like the training set size and correlation coefficients (R2, Q2) and deviation measures of the predicted values to the observed values (RMSE, MAE).
On the interactive plot training results of single structures can be inspected. E.g. the calculated descriptor values and the predicted value.

This important dialog is explained in more detail in a chapter on its own. For now please save your model with a meaningful name.

You have successfully built a prediction model on the OCHEM platform



**16.** After saving the model, you can directly apply it to predict new compounds. Before running predictions, we will investigate the model profile page in more detail.

### 3.3 Model profile



The model profile contains all the information related to the performance of the model: **statistical parameters**, **interactive scatter plot** for the single structures, links to the **data sets** used for the model and various operations, like export of the model, application of the model to new compounds, etc.

There is information about the **model name** and **public visibility**, the **predicted property** and the **training method**, used **descriptors** and **validation method**, **training duration** and **model size.**

Typical model statistics are shown in a summary table **(5)**. There are the data sets (training set and given external validation sets) with direct links to their profiles or the records respectively showing the set size. For each set the coefficient of determination „**R$^2$**", cross-validated R$^2$ called „**Q$^2$**", root mean squared error „**RMSE**" and mean absolute error „**MAE**".



Furthermore with the interactive scatterplot, showing predicted versus measured values, single records can be further investigated, e.g. comparison of predicted and measured value or inspection of calculated descriptor values for this structure.

Each point on the scatter plot is clickable and will open the "model point profile" containing the details of the respective compound from the training or the validation set.

This is a powerful feature that allows you to investigate outliers "under microscope". What are the prediction values, molecular descriptor values, the respective publication, the user who introduced this record? You can track this individually for each compound.

## Model profile ⓘ
Statistical parameters, tables, charts - all the information related to the model.

**Overview** | Applicability domain

Model name: Ames levenberg , published in Applicability domain for <I>in silico</I> models to achieve accuracy of experimental measurements
Public ID is 1

Predicted property: **AMES**
Training method: ASNN

[OEstate]
Correl. limit: 0.95 Variance threshold: 0.0,
Maximum value: 999999,
[AMES with weight 1.0
(classes weights: [inactive*0.5, active*0.5])
]Levenberg, 1000 iterations, 3 neurons
ensemble=100 k=0 additional param
PARALLEL=10
5-fold cross-validation
-
79 pre-filtered descriptors
Levenberg, 1000 iterations, 3 neurons
ensemble=100 k=16

*Calculated in 2614 seconds*
*Size: 1019 Kb*

| Data Set | # | Accuracy | Balanced accuracy | MCC | AUC |
|---|---|---|---|---|---|
| ○ Training set: Ames challenge training | 4359 records | 77.7% ± 0.6 | 77.5% ± 0.6 | 0.55 ± 0.01 | 0.854 ± 0.01 |
| ○ Test set: Ames challenge test [x] | 2181 records | 79.6% ± 0.8 | 79.5% ± 0.9 | 0.59 ± 0.02 | 0.875 ± 0.01 |

Show ROC curves

| Real↓/Predicted→ | inactive | active | Hit rate |
|---|---|---|---|
| inactive | 1512 | 504 | 0.75 |
| active | 467 | 1876 | 0.801 |
| Precision | 0.764 | 0.788 | |
| Training (Original) | | | |

| Real↓/Predicted→ | inactive | active | Hit rate |
|---|---|---|---|
| inactive | 789 | 220 | 0.78 |
| active | 225 | 947 | 0.81 |
| Precision | 0.78 | 0.81 | |
| Test (Original) | | | |

Number of compounds ignored because of errors in original model = 2

📄 Download model statistics    🖥 View configuration XML    🖥 Export configuration XML

**APPLY THE MODEL TO NEW COMPOUNDS**

For classification models, the model profile shows different statistical parameters. These are:

- accuracy, balanced accuracy, MCC, AUC,
- ROC curves
- Confusion matrices, where you can see the number of false positives, false negatives and so on.

## 3.4 Applicability domain



Each prediction in OCHEM is complemented with an accuracy estimate.

The key concept used for the accuracy estimation is so called **distance to model** (DM). DM is any measure of prediction uncertainty correlated with the prediction accuracy.

Usually, the prediction accuracy falls as DM grows, which is shown on so called **accuracy averaging plots** (shown on the left). The accuracy can be averaged in several manners: via bin-based averaging (used for regression models) and sliding window averaging (used for classification models).



With the Applicability Domain (AD) user interface different display options of the AD can be selected. There is the "Distance to model" type, the "averaging type", sliding "window size" and the label show at the x-axis.

## Smart outlier detection using AD plots



A practical feature of the applicability plots with their estimated distance to model is the automatic outlier selection according to a significance range of the p-value.

Since one measure of the distance to model follows a normal distribution, certain records can be selected and excluded from the training set.

Sometimes these "outliers" are not easily detectable as such just by inspecting the scatterplot.

# OCHEM tutorial handout

## Model export



1. It is possible to export the data related to your model by clicking "Download model statistics in Excel format". The appearing dialog allows you to select detailed info for the training and validation set – the molecular structures, identifiers, predicted and measured values, prediction accuracies, etc.

2. You can export this data in Excel, CSV, SDF or R formats. For this tutorial, please try to export an Excel file.

# 4 Working with data

In this chapter, we learn how to upload data using

- Batch data upload
- Advanced basket management

## 4.1 Batch data upload



Although you can introduce each record individually, this is usually not practically feasible. Instead, it is convenient to upload hundreds of thousands of records from external files, e.g. Excel or SD files. This can be done using the "Batch data upload" utility.

In this tutorial, we are going to upload about a thousand records for aquatic toxicity (namely, growth inhibition concentration for T. pyriformis).

**1.** Select the "Batch data upload" item in the "Database" submenu of the main OCHEM menu. You will open the first page of the "Batch upload wizard".



**2.** Select your provided SD-file in the "Upload file" field. The tool supports SDF and XLS file formats. To remind you, this file contains about 1,000 measured values for the growth inhibition assay.

**3.** Make sure you select "make the uploaded records hidden" to avoid data conflicts with the other course participants.

**4.** Hit "Upload" to continue.

**5.** The second page of the wizard is the file review page with "column remapping" tool. Here you can preview the first few lines of your uploaded file and see which columns were recognized by the system. On this page, you also have the possibility to reassign column names and select/deselect columns for upload.

**Note:**

Column headers are colour coded. Green means recognized by the system, red means not recognized. The property is dark green if it is already in the system.

Columns can be remapped by clicking on the column header



**6.** For example, the column holding the data values is named "UNIT …" in the uploaded file. We need to specify that these values represent the unit for the "Aqueous toxicity" property. Click on the red unrecognized "Unit …" column header and select "Known column" from the popup menu. Then select unit from list of known columns.

Batch Upload 3.0 - File preview and column remapping
Preview your data, select the sheet and the columns you would like to upload

T_pyriformis_all.sdf

| ☑ MOLECULE | ☑ SMILES | ☑ N | ☑ log(IGC50-1) | ☑ unit | ☑ CASRN | ☑ NAME | ☑ COMME |
|---|---|---|---|---|---|---|---|
| 1 2 3 11 11 0 0 0 0 9... | CCC(O)CC1=CC=CC=C1 | 1 | -0.16 | -log(mmol/L) | 120055-09-6 | | |
| 1 2 3 13 13 0 0 0 0 9... | OCCCCCCC1=CC=CC=C1 | 2 | 0.87 | -log(mmol/L) | 2430-16-2 | | |
| 1 2 3 11 11 0 0 0 0 9... | OCCCCC1=CC=CC=C1 | 3 | 0.12 | -log(mmol/L) | | 4-Phenyl-1-butanol | |
| 1 2 3 11 11 0 0 0 0 9... | CCC(C)(O)C1=CC=CC=C1 | 4 | 0.06 | -log(mmol/L) | 1565-75-9 | | |
| 1 2 3 13 13 0 0 0 0 9... | CCCCCOC1=CC=C(N)C=C1 | 5 | 0.97 | -log(mmol/L) | 39905-50-5 | | |
| 1 2 3 14 14 0 0 0 0 9... | CCCCCCOC1=CC=C(N)C=C1 | 6 | 1.38 | -log(mmol/L) | 39905-57-2 | | |
| 1 2 3 10 10 0 0 0 0 9... | CC(C)C1=CC=C(N)C=C1 | 7 | 0.22 | -log(mmol/L) | 99-88-7 | | |
| 1 2 3 11 11 0 0 0 0 9... | CCCCC1=CC=C(N)C=C1 | 8 | 1.07 | -log(mmol/L) | 104-13-2 | | |
| 1 2 3 9 9 0 0 0 0 9... | BrCCC1=CC=CC=C1 | 9 | 0.42 | -log(mmol/L) | | (2-BROMOETHYL)BENZENE | |
| 1 2 3 8 8 0 0 0 0 9... | CC1=CC=CC=C1N | 10 | -0.16 | -log(mmol/L) | | 2-methylaniline | |

**Several MOLECULE or NAME columns are present.**
**The ARTICLE column is missing, the stub unpublished article will be assigned by default**

Green titles indicate recognized columns, red titles indicate errors. Please click on the red columns and select whether the column indicates a property, condition or another column type like name, value or molecule, then select the matching entity and confirm your selection by clicking on the green button on the left.
If you have irrelevant columns in your sheet, you can leave them red and they will be ignored in the further process. If you need help, feel free to drop us an e-mail at info@eadmet.com .

Upload this sheet

---

Batch Upload 3.0 - Entity remapping
Review and remap the properties, conditions, units, articles and baskets involved in the data upload

**Database entities remapping**

Property: log(IGC50-1)

Values
Unit: -log(mmol/L), min value: -2.6656, max value: 3.34

Article: unpublished

Molecule set: default

submit

---

**7.** Note that the column header has changed from red to green (recognized unit), the header name is now just unit, and the checkbox in the column header is checked, indicating that the column will be processed by the tool.

**8.** Click the "Upload this sheet" button to proceed to page three of the wizard.

**9.** The third page of the wizard is the "entity remapping" page. You can review and change some aspects of the uploaded data (property, unit used for data upload, article, etc.)

**10.** Since no article has been specified in the data sheet, a stub "unpublished" was put instead of the article. Final corrections can be done here, e.g. correcting the unit, selection of a certain article or renaming of the basket.

**11.** With "submit" the data can be uploaded to the database. In this case the data will be introduced by default as hidden data and is only visible to the current user (recommended option for the tutorial exercise).

**Note**:
    To upload data originally published in an article click on the "Unpublished" link in the "Article" section of the page.

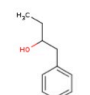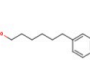Depending on the size of the uploaded set, the process may take from seconds to hours for completion (e.g. more than 50000 data points).

**12.** The fourth page of the wizard is the data preview browser. Here you can review your records and determine any errors in the data upload process.

The page holds information on the total number of records to be uploaded, the number of valid, and erroneous or duplicated records among them. You can select or deselect individual records from the upload.



**13.** Since all records being uploaded are valid, continue the upload by clicking the big "Proceed with upload" button.

The upload itself is the slowest part in the process. It may take from seconds (for a hundred records) to several hours (for a large dataset of tens of thousands of records).



**14.** The final page of the batch data upload wizard gives some statistics about the uploaded data. You can review the uploaded data in the "Experimental property browser" or download a detailed report.

**Note:**
For your convenience, the uploaded data are automatically put into a newly created basket.

## 4.2    Advanced data management



OCHEM allows combining experimental records into reusable sets. Such sets of records are referred to as **baskets**. Baskets have names and each user can virtually have as many baskets as required. Baskets are typically used as training or validation sets for the development of QSAR models.

**Basket browser and basket profile**

**1.**  Review a list of all your baskets in the basket browser accessible from *Database > Baskets* menu
The list should contain the basket created during the batch upload process from the previous tutorial.

**2.**  To open the profile of a particular basket, click on its name.

**3.**  The profile shows you brief information on the basket size, its content, articles, properties and tags. Here, you can also rename your basket or perform a number of advanced operations on it.

Please give your basket a name of your choice (e.g., "T. Pyriformis tutorial dataset") and click save.

**Splitting the basket**

**4.**  Click the link to randomly split your basket into two subsets.

**Basket splitter**

You are going to split the basket **T Pyriformis tutorial dataset** into two new baskets.
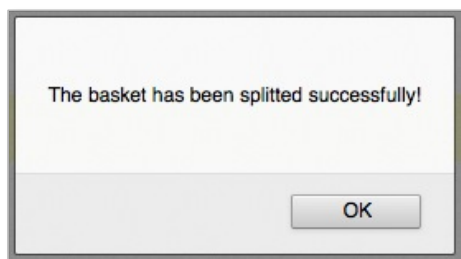
**Provide the basket names**

Basket 1: T Pyriformis tutorial dataset (training)
Basket 2: T Pyriformis tutorial dataset (test)

**Select the splitting method**

⦿ Random splitting

  Size of the validation set, in percentages: 40  %

○ Y-based splitting (not implemented yet)

Your original basket will be preserved.

( Split the basket )

**5.** Enter names for the two new baskets to be created

**6.** Select the percentage by which the molecules are divided randomly in the training and validation set.

**7.** Click "split the basket" and wait for a couple of seconds.

---

The basket has been splitted successfully!

OK

**8.** After notification that the new baskets have been created, you are forwarded to the basket browser.
**9.** The process is complete and you can see two new baskets in the basket browser.

We will use these baskets further in the model development tutorial.

---

**Basket browser** ⓘ
Browse, Compare or Join molecule set

Filter by name: [            ]  [Create new 🗋]  ☐Show public sets

1 - 15 of 20                                    15 ▾ items on page  1  of 2  >   >>

| | | |
|---|---|---|
| ☐ 📄📝🗑 Records in trash | 0 records | |
| ☐ 📄📝🗑 Selected records | 644 records | 5 pending models 🗔 |
| ☐ 📄📝🗑 T Pyriformis tutorial dataset (test) | 437 records | |
| ☐ 📄📝🗑 T Pyriformis tutorial dataset (training) | 656 records | |
| ☐ 📄📝🗑 T Pyriformis tutorial dataset | 1093 records | |

# 5. MMP Analysis

In the previous chapters, you have learned to perform the most basic steps of a QSAR modelling. In this chapter, we will review advanced features:

Interpretation of models using Molecular Matched Pairs

Matched Molecular Pair analysis to interpret models.

**1.** Click on the highlighted region to open the MMP plot



**2.** Points near to the diagonal corresponds to the MMPs which were learnt. Points in 2nd and 4th quadrants correspond MMPs for which even the sign was incorrectly predicted.

**3.** Clicking on each point will show the transformation responsible for each MMP. In the shown case the point in the 4th quadrant does not agree with qualitative changed of property (logP) for other transformations. It is likely to be an error, indeed addition of two carbon atoms CC increases logP – this is not an expected behavior.

**4.** The same approach can be used to interpret chemical reactions. Here the two groups of pairs of reactions correspond to two types of reactios.

## Acknowledgement

I thank Dr. Sushko, Dr. Novotarskyi, Mr. Körner and other members of my team for their work on the development of the OCHEM platform and preparation of the earlier version of this tutorial.