**Supplementary Information to article:**

**The development of models to predict melting and pyrolysis point data associated with several hundreds thousands compounds mined from patents**

Igor V. Tetko,[1,*] Daniel Lowe[2] and Antony J. Williams[3]

1) Helmholtz Zentrum München für Gesundheit und Umwelt (HMGU), Institute of Structural Biology, Ingolstaedter Landstrasse 1, b. 60w, D-85764 Neuherberg, Germany and BigChem GmbH, D-85764 Neuherberg, Germany i.tetko@helmholtz-muenchen.de
2) NextMove Software Limited, Innovation Centre (Unit 23), Cambridge Science Park, Cambridge CB4 0EY, UK, daniel@nextmovesoftware.com
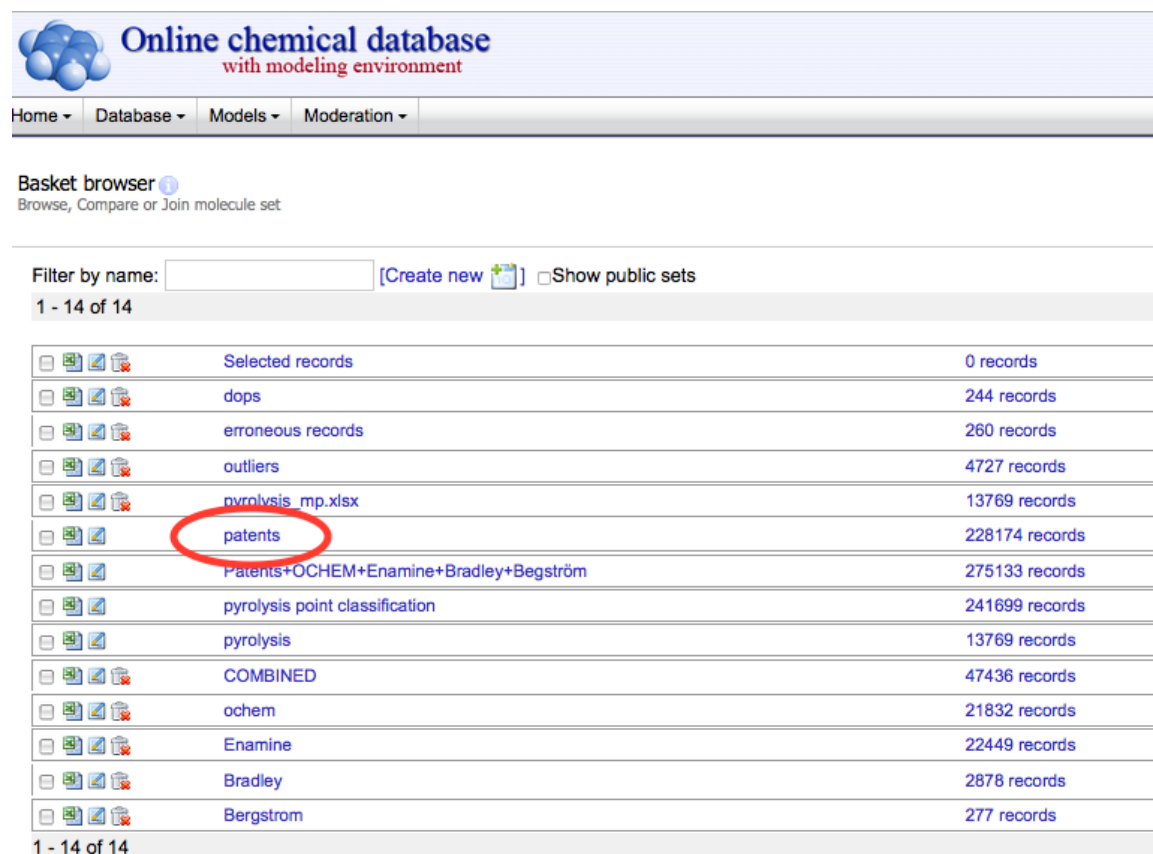3) ChemConnector Inc., 904 Tamaras Circle, Wake Forest, NC-27587, USA, Williams.Antony@epa.gov

## Description of the protocols used to develop the Melting Point consensus model.

## Overview

Below, we summarize steps used for the development of the consensus model. The illustrative Figures are followed by a brief description of the respective step.



## Step 1. Data upload.

The data for the Melting Point (MP) model are uploaded in the OCHEM database (for these steps see tutorial at https://www.youtube.com/watch?v=vdJWiS4wSaQ) and formed the basket "patents".

## Step 2. Identification of LibSVM parameters.

A preliminary analysis using a grid search (see Methods) was run to select the optimal parameters for the SVM model. It included the following several substeps.



## Step 2.1 Starting model development.

## Online chemical database
### with modeling environment

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

### Create a model ⓘ
Select the training and validation sets, the machine learning method and the validation protocol

**Select the training and validation sets:**

Training set *(required)*: patents [details]
Add a validation set

The model will predict this property:
Melting Point using unit: [ °C ◇ ]

**Choose the learning method:** ⓘ

*Suggested modeling methods:*
- ○ ASNN (ASsociative Neural Networks)
- ○ FSMLR (Fast Stagewise Multiple Linear Regression)
- ○ KNN (K-Nearest Neighbors)
- ○ Library model (A local bias correction model based on another ASNN model)
- ● LibSVM with grid-search parameter optimisation
- ○ MLR (Multiple Linear Regression)
- ○ PLS (Partial Least Square)
- ○ WEKA-J48 (C4.5 decision tree) classification only, use with bagging only
- ○ WEKA-NB (Naive Bayes) classification only
- ○ WEKA-RF (Random Forest) classification only

*Methods under development:*
- ○ Consensus model

**Model validation**
Validation method: [ No validation ◇ ]

Development of a model without validation may lead to incorrect estimation of the prediction accuracy of your model

You can create a model from template: import an XML model template or use another model as a template

[ Next>> ]

## 2.2 Selection of the dataset, method and validation protocol to identify the best set of LibSVM parameters.

Since the calculations in this step required very large CPU resources and LibSVM parameters were optimized for the whole set, "no validation" method was used. This method is normally not recommended since it can result in data overfitting. In this case only three LibSVM parameters were optimized for N = 228k compounds and thus there was no problem with overfitting. Notice that these calculations required more than 600 CPU-days of calculations.

4

## 2.3 Selection of preprocessing options.

The options included use of records with intervals and ranges for model development.

**Online chemical database**
with modeling environment

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

**Model creator**
Select model template and training set

**Select the molecular descriptors** ⓘ

**Recommended descriptor types**

- ☐ E-state
- ☐ ALogPS (2)
- ☐ GSFragment (1138)
- ☐ Dragon v. 6.0 (4885/3D)
- ☐ ISIDA fragments
- ☐ ADRIANA.Code *(211/3D)*
- ☐ CDK descriptors (246/3D)
- ☐ 'Inductive' descriptors *(54/3D)*
- ☐ MERA descriptors *(529/3D)*
- ☐ MERSY descriptors *(42/3D)*
- ☐ Chemaxon descriptors *(499/3D)*
- ☐ QNPR
- ☐ Spectrophores (144/3D)
- ☑ Structural alerts (ToxAlerts)

  **Select alerts:**

  | | |
  |---|---|
  | Publication | All articles ⌄ |
  | Endpoint | Extended Functional groups ⌄ |
  | ☐ Only approved alerts | |

**Special descriptors (scaffolds, fingerprints):**

- ☐ SIRMS
- ☐ Scaffold Hunter Descriptors
- ☐ Chemaxon Scaffolds
- ☐ Silicos-It Scaffolds
- ☐ ECFP Fingerprints
- ☐ MolPrint Fingerprints

☑ Forbid NaN and Infinite descriptor values

[ <<Back ] [ Next>> ]

# 2.4 Selection of Extended Functional Groups (EFGs) as descriptors.

**Online chemical database**
with modeling environment

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

Model creator
Select model template and training set

## Select filters of descriptors

☑ Eliminate descriptors with less than [2] unique values

☑ Delete descriptors that have absolute values larger than [999999]

☑ Delete descriptors that have variance smaller than [0.01]

☑ Group descriptors, that have pair-wise correlations Pearson's correlation coefficient $R$ larger than [0.95]

☐ Use Unsupervised Forward Selection to delete variables using the above value of multiple correlation coefficient $R$

☐ After filtering, I want to select necessary descriptors myself (advanced)

**Normalisation parameters**

| | |
|---|---|
| Descriptors normalization | To range [0, 1] (e.g., for LibSVM) |
| Values normalization | Do not normalize |

[<<Back] [Next>>]

## 2.5 The descriptor unsupervised filtering options

The normalization to [0,1] interval was selected. Other pre-processing options were used as default values.

**Online chemical database**
with modeling environment

Home ▾  Database ▾  Models ▾  Moderation ▾

Model creator
Select model template and training set

**Configure LibSVM method**

**SVM Method options**

SVM algorithm  Classic (C-SVC, epsilon-SVR)

Kernel type  Radial Basis Function

☐ **Use class weighting for classification tasks**
☑ **Enable grid search**

*Parameters are assigned values 2^(current_step) where current_step = (min, min+step, ..., max)*

Cost min, max, step:    -10  , 10  , 2
Gamma min,max,step:   -10  , 10  , 2
Epsilon min,max,step:  -16  , 10  , 2
Grid search set size
(fraction of training set):  1
PARALLEL:   1300

☐ **Configure advanced options**

<<Back    Next>>

## 2.6 Grid search to detect optimized parameters.

The calculations were configured to run in parallel on 1300 servers (option PARALLEL) and using the whole training set (fraction = 1).

## Online chemical database
### with modeling environment

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

## Model creator
Select model template and training set

### Start calculation of the model

Now we are ready to start calculation.

Please provide the name for your model: Melting Point_LibSVM_[StructuralAlerts], 248469

☑ Save models

Task priority:
○ High priority (please, use for fast tasks only)
● Normal priority
○ Low priority (for long tasks)

[ <<Back ] [ Start calculation>> ] [ Discard ]

## 2.7 The calculations are launched.

## 2.8 Selection of Pending task window to monitor the execution of tasks.



## 2.9 Monitoring of the calculations in the PendingTask browser.

## 2.10 The model is calculated.

The optimal parameters for the training set were calculated. They could be exported using XML icon at the left side of model row.



## 2.11 The model is open (green button on the previous plot) and saved.

The optimized parameters for LibSVM are shown. They are part of the model template, which can be exported using the "Export configuration XML" button. The exported template can be used to develop new models (see step 2.1 – "import an XML model template") using the same setting.

Online chemical database
with modeling environment

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

Create a model
Apply a model
Create multiple models
Open predictor
Upload a linear model
Upload a stub model

View pending tasks
View published tasks

SetCompare utility
MolOptimiser
Calculate descriptors
Descriptors storage

Comprehens
Create multiple

The comprehensive ... u to simultaneously run multiple models with different machine learning methods, molecular descriptors and validation protocols.
Please note that ru... equire significant computational resources and time.

**Select the traini...**

Training set *(requ...*
Add a validation ...

Select the methods you want to use for the modeling:

| Method | Descriptors | Descriptor selection | Model validation |
|---|---|---|---|
| [all] [none] | [all] [none] | [all] [none] | [all] [none] |
| ☐ ANN | ☑ CDK | ☐ Unsupervised forward selection | ☑ 5-fold cross-validation |
| ☑ ASNN (bias correction) | ☐ Dragon v.6 (all blocks) | ☑ Pairwise decorrelation (r < 0.95) | ☐ 5-fold cross-validation (stratified - classification only) |
| ☐ KNN | ☑ OEstate and ALogPS | ☐ No UFS [edit] [x] | ☐ Bagging with 64 models |
| ☐ LibSVM | ☐ ISIDA Fragments (Length 2 - 4) | | ☐ Bagging with 64 models (stratified - classification only) |
| ☑ FSMLR | ☐ GSFrag | +add a custom template | |
| ☐ MLRA | ☐ Mera and Mersy | | +add a custom template |
| ☑ PLS | ☐ Chemaxon descriptors | | |
| ☐ WEKA-RF (classification only) | ☐ Inductive Descriptors | | |
| ☐ WEKA-J48 (classification only) | ☐ Adriana | | |
| ☐ LibSVM parallel [edit] [x] | ☐ Spectrophores | | |
| ☐ LibSVM optimised 2 [edit] [x] | ☐ QNPR (SMILES - length 1 - 3 threshold 5) | | |
| | ☐ OEstate [edit] [x] | | |
| +add a custom template | ☐ Drafon 5.4 [edit] [x] | | |
| | ☐ MolPrint [edit] [x] | | |
| | ☐ Functional groups | | |
| | ☐ ECFP4 | | |
| | ☑ Mol. weight + # of carbons: baseline model | | |
| | ☐ OEstate counts [edit] [x] | | |
| | ☐ SIRMS | | |
| | +add a custom template | | |

Show advanced options>>

Considering the selection above, **9 models** will be created.

Create the models

## 2.12 The optimal parameters of the model are used to create a new method template.

The "Comprehensive Modeling" mode was opened and a template was created using the "add a custom template on "Models/Create multiple models" menu. The previously selected model with optimized parameters was used as the template. An experienced user can also manually edit the XML file in the browser. In this example template with optimized LibSVM parameters was created. In a similar way new templates for descriptors, unsupervised descriptor selection and model validation can be created. They could be combined to create multiple models.

## 3. The development of multiple models is started using model template with optimized LibSVM parameters.

For this article we used default settings for descriptors, descriptor selection and validation protocols as shown in the Figure above.

Models calculation was monitored in the Pending Task menu (see 2.8 and 2.9). Once models were finished, they were stored (each model individually, see step 2.11). After that they were used to develop a Consensus model using the "Consensus model" template and the same steps (2.1 and 2.2) used for the LibSVM model.

## 3.1 Consensus model developed using the "Simple average" option.



## 3.2 Selection of saved models for averaging in "Consensus model".

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

**Database menu:**
- **Compound properties**
- Properties
- Conditions
- Units
- Articles/Books
- Journals
- ToxAlerts ▸
- MatchedPairs ▸
- **Baskets**
- Tags
- Set area of interest...
- User-related changes
- **My data exports**
- **Batch data upload**
- Trash

Basket
Browse, 

Filter
1 - 14

[Create new] ☐ Show public sets

| | records | 0 records | |
| | | 244 records | 1 models (+13 pending) |
| | records | 260 records | |
| | | 4727 records | |
| | .xlsx | 13769 records | |
| | | 228174 records | 36 models (+1 pending) [overview] |
| | HEM+Enamine+Bradley+Begström | 275133 records | 8 models (+5 pending) [overview] |
| | pyrolysis point classification | 241699 records | 7 models (+2 pending) [overview] |
| | pyrolysis | 13769 records | 12 models [overview] |
| | COMBINED | 47436 records | |
| | ochem | 21832 records | |
| | Enamine | 22449 records | |
| | Bradley | 2878 records | |
| | Bergstrom | 277 records | |

1 - 14 of 14

## Step 4. The calculated models are accessed in the table – like Comprehensive Modeling (CM) view from "Baskets" web page.

Home ▾ | Database ▾ | Models ▾ | Moderation ▾

Molecule sets X | Basket models summary X | Basket models summary X | Basket models summary X | Basket models summary X | Model profile X

Multiple models overview

Predicted property: Melting Point
Training set: patents (4 different versions detected) ⓘ

Metrics [RMSE - Root Mean Square Error ▾] for [Training + Excluded ▾] Validation: [All validation protocols ▾]

| | LibSVM | LibSVM (tr. set. 2) | LibSVM (tr. set. 3) | LibSVM (tr. set. 4) | MLRA |
|---|---|---|---|---|---|
| CDK (constitutional, topological, geometrical, electronic, hybrid) | 38.93 | 38.87 | 38.84 | 38.88 | 45.3 |
| Fragmentor (Length 2 - 4) | 38.49 | 38.22 | 38.33 | 38.43 | 45.17 |
| ChemaxonDescriptors (7.4) | 40.07 | 40.13 | 40.1 | 40 | 48.33 |
| QNPR (SMILES - length 1 - 3 threshold 5) | 39.74 | 39.24 | 39.44 | 39.69 | 46.4 |
| OEstate | 38.26 | 38.04 | 38.1 | 38.13 | 46.35 |

### Step 4.1 An overview of the developed models in the Comprehensive Modeling (CM) window.

Models developed with individual sets of descriptors are shown. The columns correspond to sets formed by the exclusion of outlying molecules with different p-values. The models developed using MLRA calculate with much lower accuracy.