

## Problem Definition and Goals

**Problem:** The performance of QSAR models is hampered by the *quality* of the underlying data.

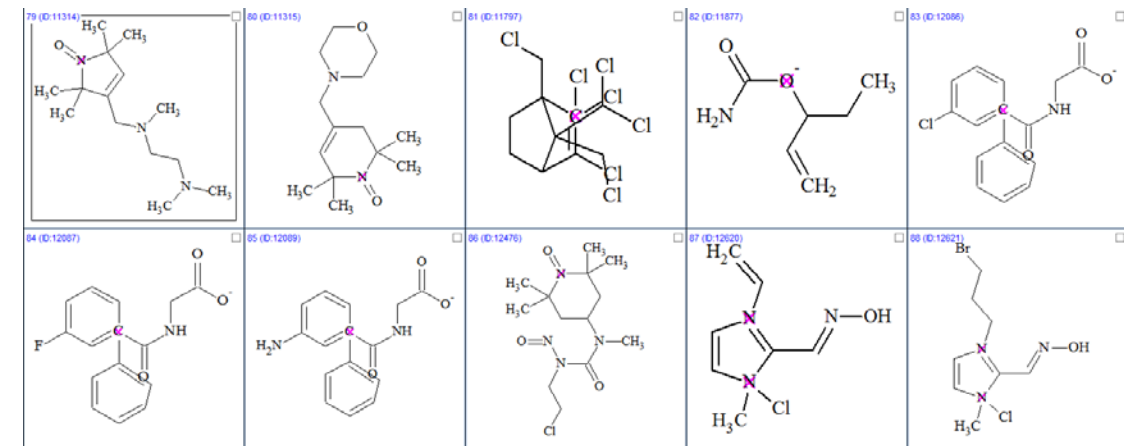
**Goals:** To examine the quality of the data underlying the EPI Suite prediction models using both manual and automated methods. To use the curated data to develop new prediction models and examine the influence of data quality and quantity, on the resulting QSAR predictions. To make the data and models publically available.

## Abstract

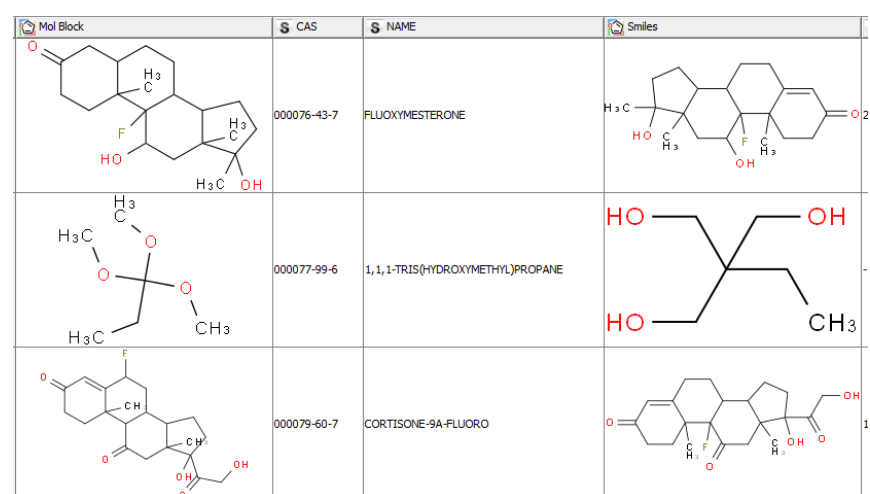
The accuracy of QSAR models is critically dependent on the quality of available data. As part of our efforts to develop public platforms to provide access to predictive models, we have attempted to discriminate the influence of the quality versus quantity of data available to develop and validate QSAR models. We have focused our efforts on the widely used EPI Suite<sup>1</sup> data (PHYSPROP database) that was initially developed over two decades ago. Specific examples of quality issues for PHYSPROP data include multiple records for the same chemical structure with different measured property values, inconsistency between the structure, chemical name and Chemical Abstracts Service Registry Number (CASRN) for single records, the inability to convert the SMILES strings into chemical structures, hypervalency in the chemical structures and the absence of stereochemistry for thousands of data records. Relative to the era of EPI Suite development, modern cheminformatics tools allow for more advanced capabilities in terms of chemical structure representation and storage, as well as enabling automated data validation and standardization approaches to examine data quality. This poster reviews both our manual and automated approaches to examining key datasets related to PHYSPROP data. This includes approaches to validate between chemical structure representations (e.g., Mol-Block and SMILES) and identifiers (chemical names and CASRN), as well as approaches to standardize the data into QSAR-ready formats for modeling. We have quantified and segregated the data into various quality categories to allow us to thoroughly investigate the resulting models that can be developed from these data slices and to examine to what extent efforts into the development of large high-quality datasets have the expected pay-off in terms of prediction performance.

## Source Data and Example Errors

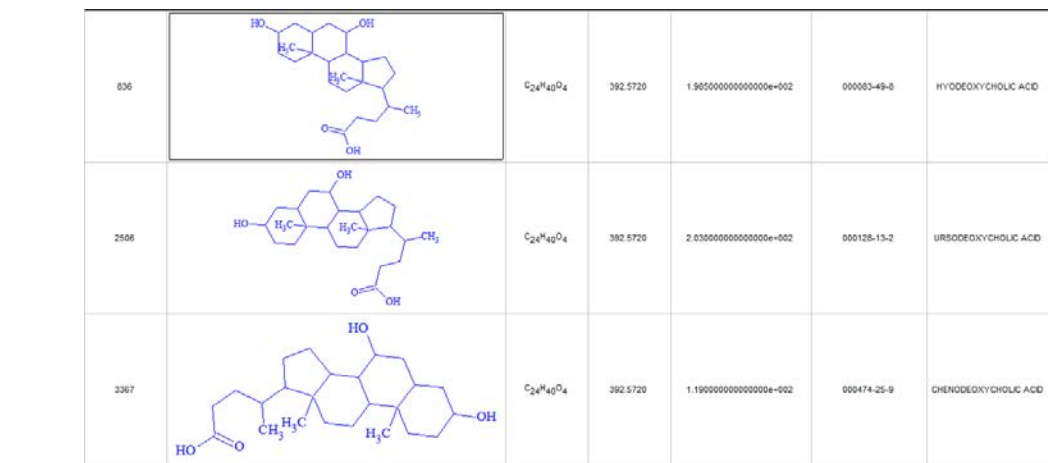
- The PHYSPROP data were sourced online<sup>2</sup> as SDF files. A total of 13 endpoints were represented, including LogKow, water solubility, melting point, boiling point, and others. The largest dataset, LogKow, contained over 15,800 individual chemicals. Each data point included the Mol-Block, SMILES, CASRN, Name, LogKow value and, where available, a reference. This dataset was chosen as representative for examining the quality of data.
- A manual examination of the data revealed a number of issues: e.g., SMILES and Mol-Block did not agree; CASRN did not match the correct structure; SMILES could not be converted; a single chemical structure would be listed multiple times with different property values. Example errors across various PHYSPROP datasets are shown below.



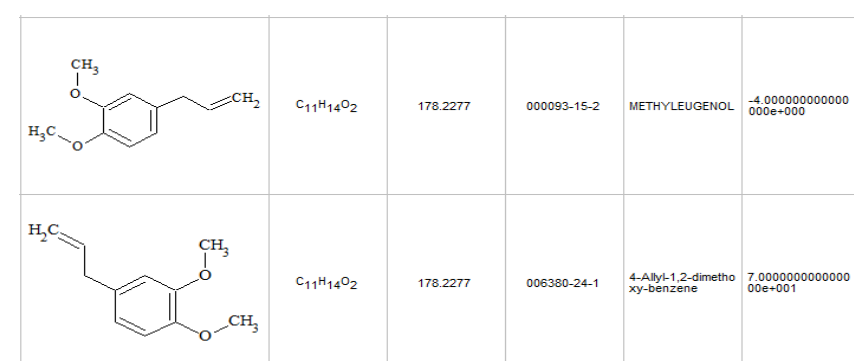
Examples of hypervalency in LogKow dataset



Different structures in Mol-Block and SMILES



Equivalent structures, different CASRN, names, and values in MP dataset

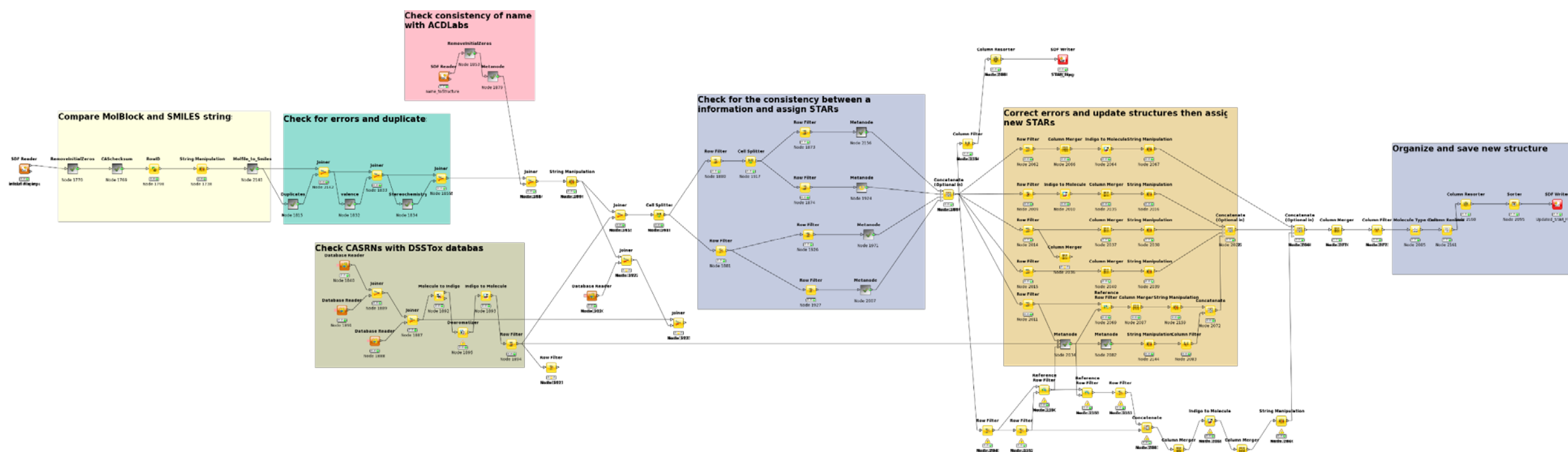


Equivalent structures but with different CASRN, names and values in MP dataset: -4 and + 70°C

## Automated Analysis Using KNIME

The manual investigation of the data allowed us to develop a KNIME<sup>3</sup> workflow for automated processing. This workflow was derived from earlier work by Mansouri *et al.*<sup>4</sup> and is represented in the figure below as a series of blocks representing, for example:

- Compare Mol-Block and SMILES (2268 different)
- Check for duplicates (657 structures, 531 names)
- Check CASRN Numbers (3646 invalid CASRN)
- Check names against dictionary (555 invalid)
- Assign Quality flags based on consistency among data fields



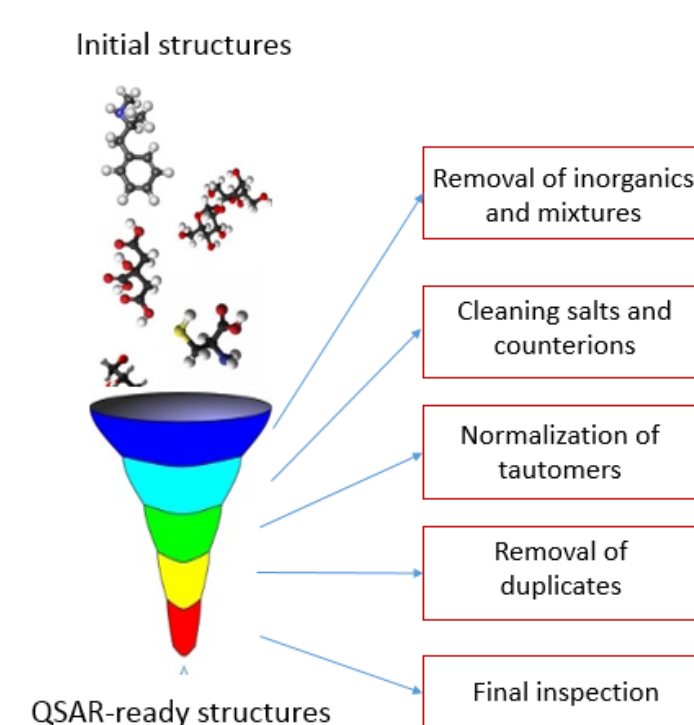
The KNIME workflow for automated processing of PHYSPROP data

The KNIME workflow was used to insert various levels of Quality Flags indicating consistency between chemical structure formats and identifiers. The consistency flag definitions and distribution are summarized below for the >15k chemicals.

4 STAR ENHANCED:	Name/CASRN/Mol/SMILES added Stereo:	550
4 STAR:	3 of 4 Name/CASRN/Mol/SMILES:	5967
3 STAR ENHANCED:	3 of 4 Name/CASRN/Mol/SMILES added Stereo:	177
3 STAR:	3 of 4 Name/CASRN/Mol/SMILES:	7910
2 STAR PLUS:	2 of 4 Name/CASRN/Mol/SMILES/Tautomer:	133
2 STAR:	2 of 4 Name/CASRN/Mol/SMILES:	1003
1 STAR:	No two fields consistent	379

Predictive models were developed using only the 3 and 4 star curated data.

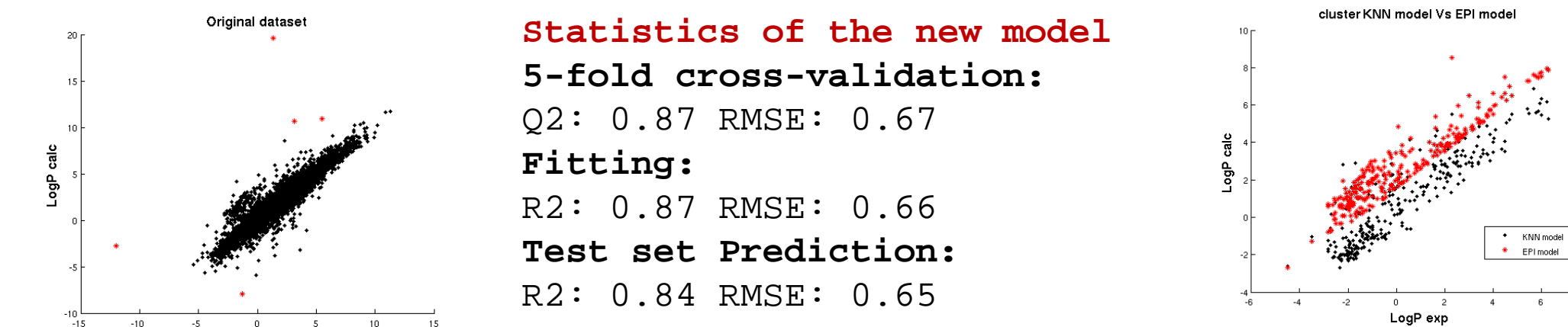
## Preparing data for QSAR Modeling



For the purposes of QSAR modeling, the 3 and 4 STAR datasets were processed through a KNIME workflow. This processing removed inorganics and mixtures, processed salts into neutral forms (except for melting point data), normalized tautomers, and removed duplicates. The resulting “QSAR-ready” file(s) were modeled using Genetic Algorithm-Partial Least Squares with 5-fold cross validation and utilizing 2D PaDEL<sup>5</sup> molecular descriptors. Multiple modeling runs (100) produced the best models using a minimum number of descriptors. The models for all 13 endpoints are available as both Windows and Linux executable binaries and as a C++ library that can be called by a separate application.

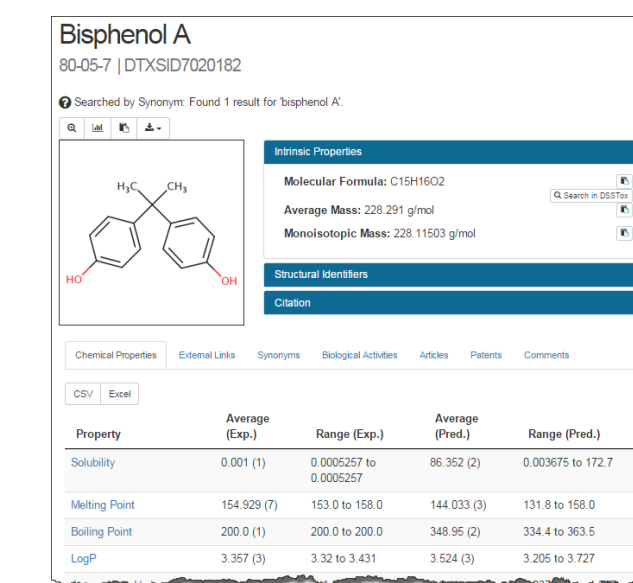
## Model Performance

The LogKow prediction model delivered by EPI Suite used a smaller dataset (of 2700 chemicals). The curation of the available data, utilization of a larger dataset (>14,000 chemicals) and application of novel machine-learning approaches produced a better and simpler model with only 10 descriptors. The figures below illustrate the difference between the original EPI Suite model and the newly derived predictive model. The red data points indicate the outliers from the original modeling approach, the majority not included in the original training set.



## Accessing ~1 Million Predicted Properties Online

The iCSS Chemistry Dashboard (ICD) is a public EPA-hosted web application providing access to over 700,000 chemicals from EPA's DSSTox database<sup>6</sup>. It integrates experimental and predicted data and is a hub to other NCCT apps and web-based resources. The 3 and 4 STAR curated experimental data are accessible via the ICD application. All chemicals were also passed through the collection of NCCT prediction models and results will be freely available at <http://comptox.epa.gov/dashboard> in April 2016.



Chemical Properties: LogP						
		Average	Range			
Experimental		3.307 (3)	3.32 to 3.431			
Predicted		3.524 (3)	3.205 to 3.727			
Property	Raw Result	Mean Result	Minimum Result	Maximum Result	Result Unit	Source
Estimated Log Kow	3.64	3.64	3.64	3.64		predicted EPI SUITE
LogP	3.431	3.431	3.431	3.431		experimental Vitas-M
LogP	3.727	3.727	3.727	3.727		predicted ACD/Labs
Measured Log Kow	3.32	3.32	3.32	3.32		experimental EPI SUITE
Octanol-water partition coefficient	3.32	3.32	3.32	3.32		experimental CURATED_PHYSPROP

## Future Work

- Release NCCT models as interactive online prediction tools in the near future via the ICD.
- Integrate the suite of EPA T.E.S.T.<sup>7</sup> physchem and toxicity prediction models to expand the collection of available models.
- Source additional data to expand the training sets underpinning the prediction algorithms.

## References

- EPI Suite: <http://www.epa.gov/tsc-screening-tools/epi-suite-tm-estimation-program-interface>
- PHYSPROP Data: [http://esc.syrres.com/interkow/EpiSuiteData\\_ISIS\\_SDF.htm](http://esc.syrres.com/interkow/EpiSuiteData_ISIS_SDF.htm)
- KNIME: <https://www.knime.org/>
- Mansouri *et al.* CERAPP: Collaborative Estrogen Receptor Activity Prediction Project, *Environ Health Perspect*; DOI:10.1289/ehp.1510267
- PaDEL descriptors, <http://padel.nus.edu.sg/software/padeldescriptor/>
- EPA Distributed Structure-Searchable Toxicity (DSSTox) Database, <http://www.epa.gov/chemical-research/distributed-structure-searchable-toxicity-dssto-database>
- EPA Toxicity Estimation Software Tool (T.E.S.T.) software <http://www.epa.gov/chemical-research/toxicity-estimation-software-tool-test>