# Establishing Best Practices for Water Solubility Dataset Curation

Charles Lowe[1], Gabriel Sinclair[1,2], Christian Ramsland[1,2], Todd Martin[1], Christopher Grulke[1], and Antony J. Williams[1]

[1]Center for Computational Toxicology and Exposure, Office of Research and Development, U.S. Environmental Protection Agency (U.S. EPA), Research Triangle Park, North Carolina, USA
[2]Oak Ridge Associated Universities (ORAU), Oak Ridge, Tennessee, USA.

ORCID: 0000-0001-9151-6157
Charles Lowe I lowe.charles@epa.gov I 919-541-5618

## Problem Definition and Goals

**Problem**: There are numerous peer-reviewed publications and public websites that contain experimental data that could be used to improve existing QSAR/QSPR models. Commonly these data are not available in an ideal form: often limited to PDF supplementary info files for publications (with names or CASRNs and no electronic structure format). However, when aggregation of these data has been attempted curation has been necessary.

**Goals**: Provide a *de facto* dataset for water solubility data that can be used to build multiple models and eventually a consensus model. Determine erroneous records through curation and validation of chemical identifiers, as well as standardize solubility values and exclude outlying values using statistical approaches and cutoff values. Make these data available as downloadable data for use in QSAR/QSPR models and reuse in other databases.

## Dataset Assembly

| Database | URL | Original No. of Records | Curated No. of Records |
|---|---|---|---|
| AqSolDB | https://doi.org/10.1038/s41597-019-0151-1 | 9559 | 5657 |
| Bradley Dataset | http://dx.doi.org/10.1021/ci800406y | 3892 | 2909 |
| eChemPortalAPI | https://echa.europa.eu/ | 7037 | 2722 |
| Ochem | https://ochem.eu/ | 27928 | 18499 |
| PubChem | https://pubchem.ncbi.nlm.nih.gov/ | 45920 | 20701 |
| OFMPub | https://ofmpub.epa.gov/oppthpv/ | 386 | 213 |
| OPERA | ftp://newftp.epa.gov/COMPTOX/Sustainable_Chemistry_Data/Chemistry_Dashboard/PHYSPROP_Analysis/ | 4839 | 4068 |
| EPISuiteISIS | http://esc.syrres.com/interkow/EpiSuiteData_ISIS_SDF.htm | 4755 | 3689 |
| LookChem | https://www.lookchem.com/ | 1820 | 907 |
| QSARDB | https://qsardb.org/ | 19 | 8 |
| Chemical Book | https://www.chemicalbook.com/ | 6 | 0 |

**Table 1:** The nine databases (and two journal articles), their corresponding URLs, and the original number of records found. The curated number of records was obtained using the workflow shown in **Figure 3**.
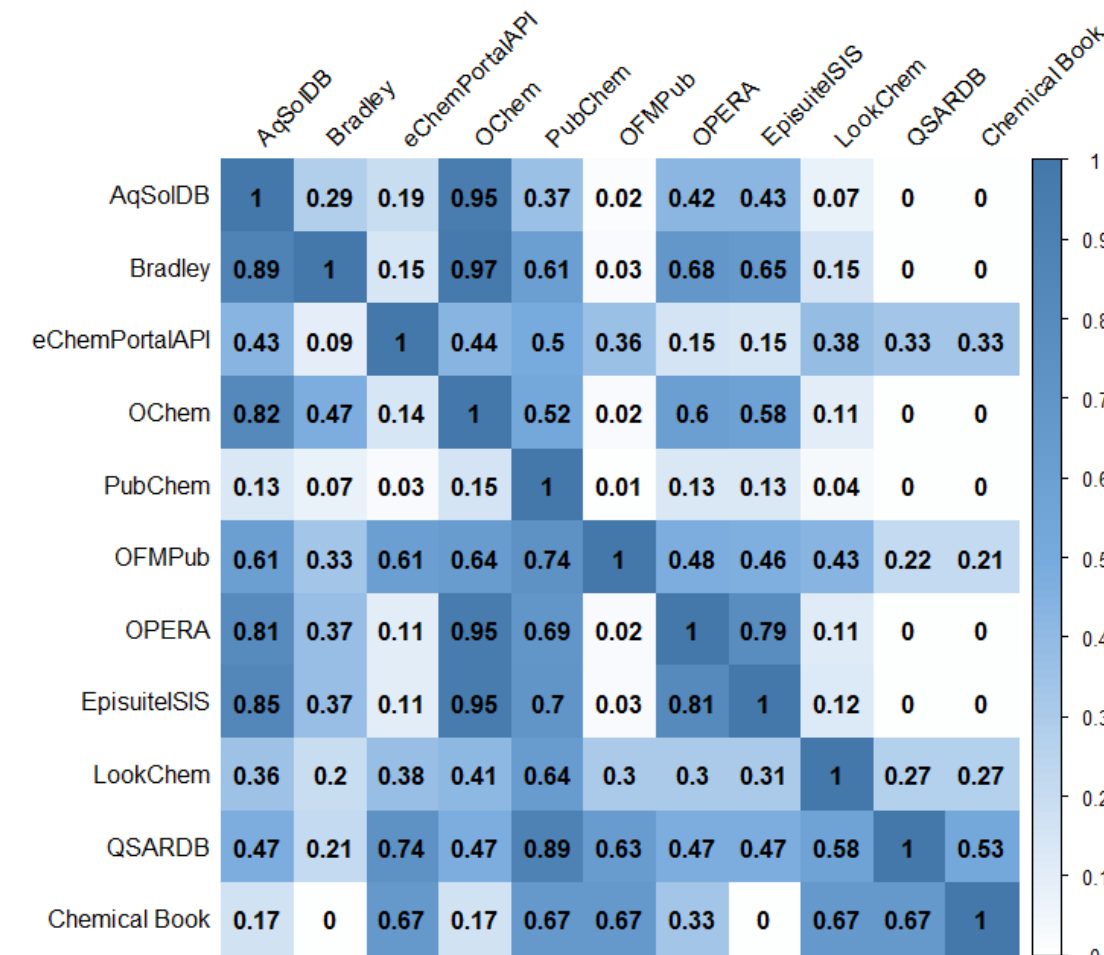


**Figure 1:** Redundancy matrix showing the intersection of chemicals between datasets as a fractional value. While some of the databases have significant overlap (i.e., OPERA with OCHEM, EPI Suite with AqSolDB), no database perfectly overlaps with another. The significant overlap between databases allows for checks of parity, where ambiguously-represented chemicals may be corrected or removed.
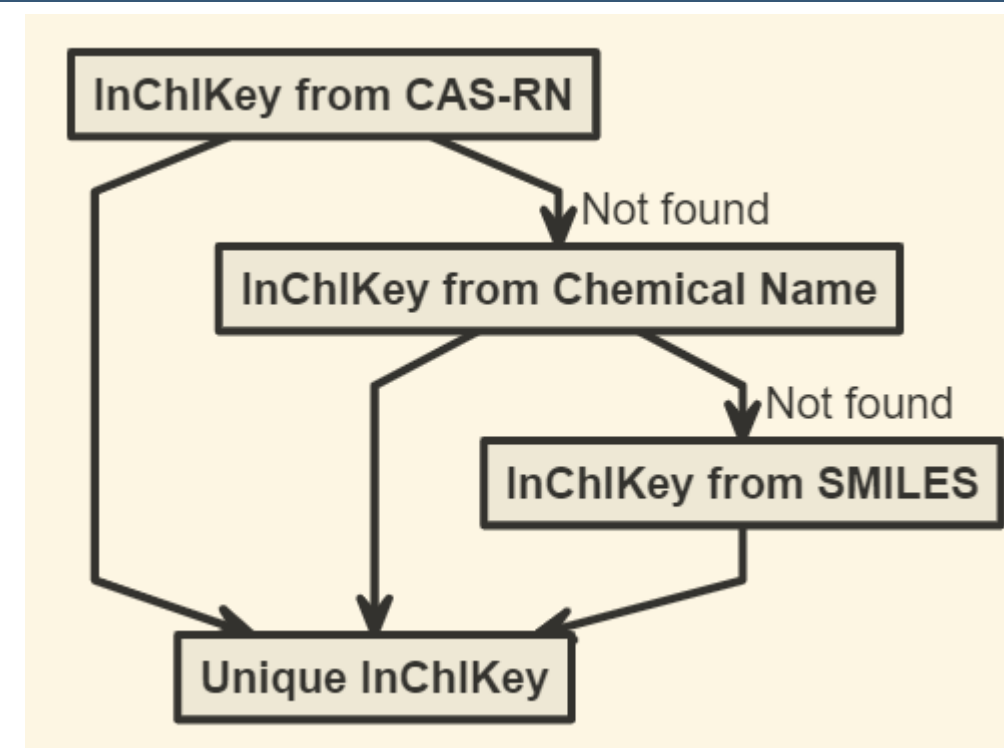
## Dataset Curation



**Figure 2:** This diagram shows the workflow for selecting a unique chemical identifier for each dataset entry. (above) InChIKeys are determined via a search of OPSIN or ACD/labs software.
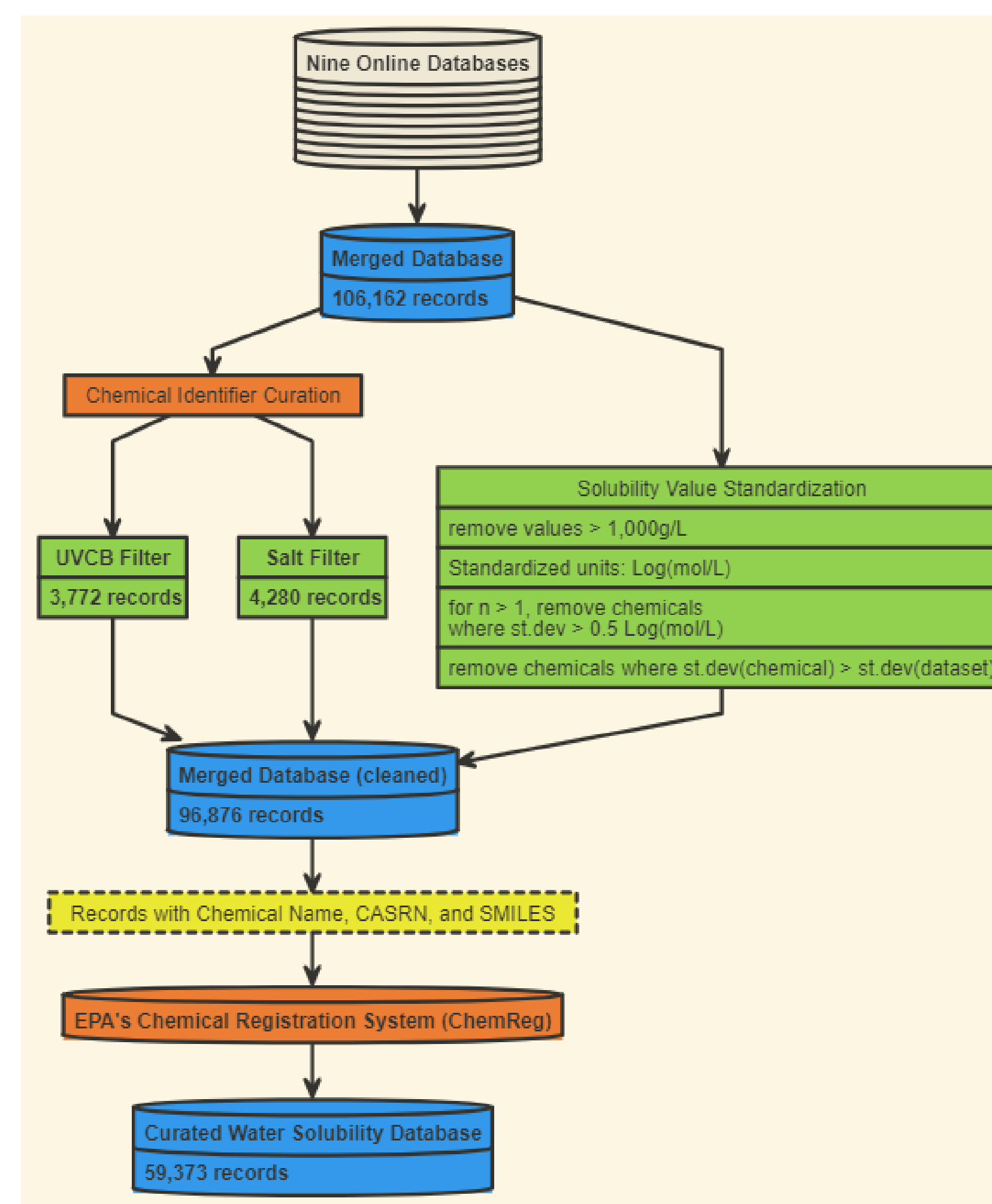


**Figure 3:** This diagram shows the workflow for curation and standardization of the dataset. Chemical identifier curation is achieved as shown in **Figure 2**.

## Main Results

- In total, 49,804 unique chemicals mapped to 47,121 QSAR-ready structures*.
- 22,675 unique chemical structures mapped to 22,365 QSAR-ready structures* in EPA's chemical registration system.
- Examples of curation issues: multiple CAS-RNs or names per record (not UVCBs), truncated chemical names, UVCBs given a single chemical structure, inverted signs (negative instead of positive) that are correctable based on database parity.

*desalted, de-isotoped, stereo-neutral forms of chemical structures

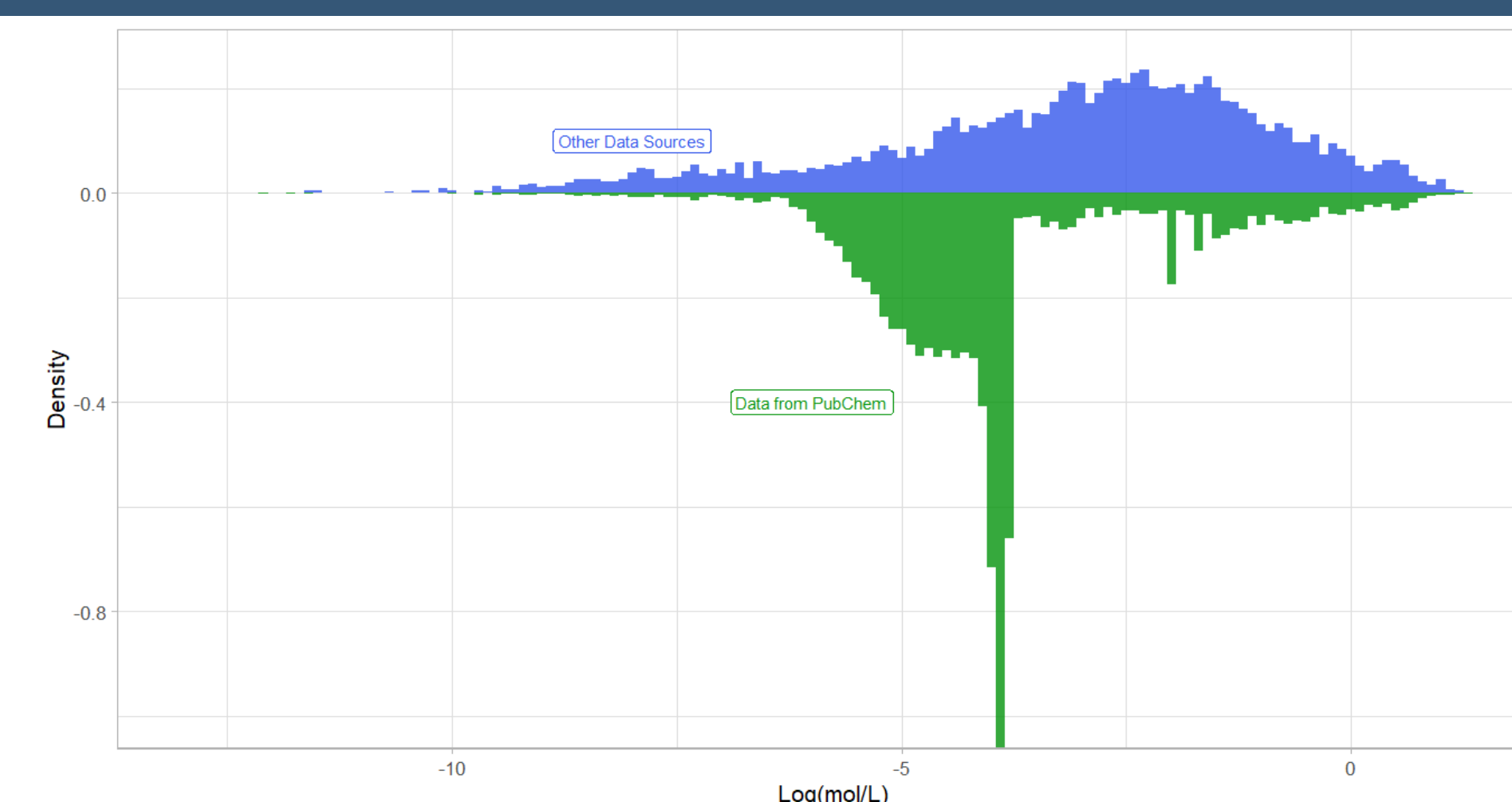## Preliminary Modeling of Curated Dataset



**Figure 4:** A density plot showing the different distributions of solubility values for the data from PubChem versus the data from other sources. The presence of numerous drugs in PubChem has imbalanced the overall dataset toward the fit-for-purpose region, [3 Log(mol/L) to 5 Log(mol/L)].
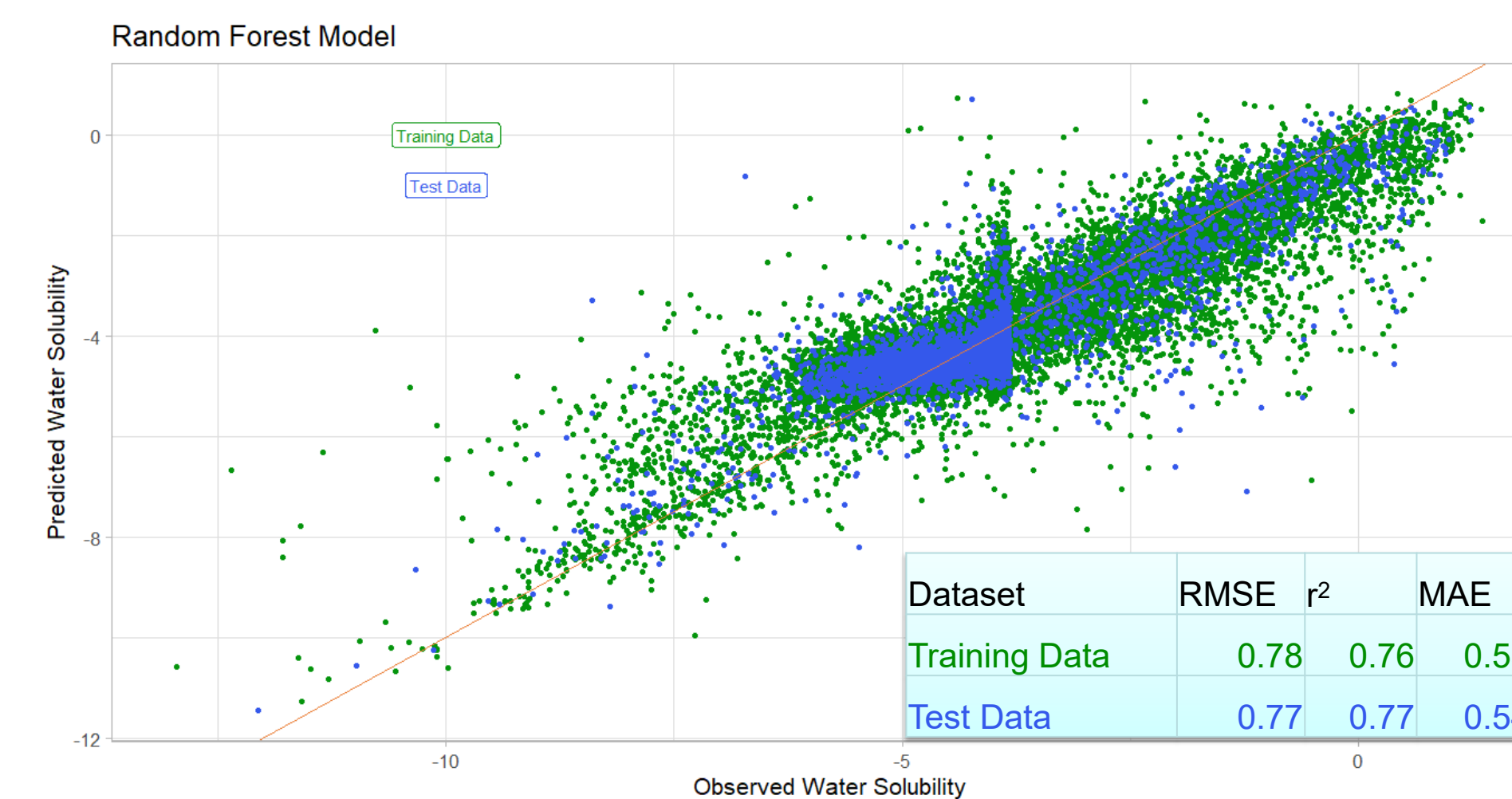


| Dataset | RMSE | $r^2$ | MAE |
|---|---|---|---|
| Training Data | 0.78 | 0.76 | 0.55 |
| Test Data | 0.77 | 0.77 | 0.54 |

**Figure 5:** This plot shows the correlation of experimental values to prediction results obtained from a random forest model constructed using 85% of the 22,365 QSAR-ready structures/solubility value pairs and PaDEL descriptors.