

# BigSolDB: Solubility Dataset of Compounds in Organic Solvents and Water in a Wide Range of Temperatures

Lev Krasnov<sup>1</sup>, Simon Mikhaylov<sup>1</sup>, Prof. Maxim V. Fedorov<sup>2</sup>, and Dr. Sergey Sosnin<sup>3,\*</sup>

<sup>1</sup>Department of Chemistry, Lomonosov Moscow State University, 119991 Moscow GSP-1, 1-3 Leninskiye Gory, Russia

<sup>2</sup>Sirius University of Science and Technology, Olimpiyskiy ave. b.1, Sochi, 354000, Russia

<sup>3</sup>Department of Pharmaceutical Sciences, Faculty of Life Sciences, University of Vienna, Josef-Holaubek-Platz 2, 1190 Vienna, Austria

\*corresponding author(s): Dr. Sergey Sosnin (sergey.sosnin@univie.ac.at)

## ABSTRACT

Solubility is crucial in organic chemistry and holds significant value in the field of medicinal chemistry. Employing computational and QSPR modeling for solubility estimation is favorable as it reduces experimental costs. However, high-quality experimental data is essential for training these QSPR models. In our study, we compiled a dataset consisting of 54,273 experimental solubility values within a temperature range of 243.15 to 403.15 K in various organic solvents and water. This dataset can be used as a reference for individual values or training solubility QSPR models. We conducted a statistical analysis and identified prevalent patterns in the data. Furthermore, we developed an interactive, parametric t-SNE-based tool to explore the chemical space of solutes. Utilizing this tool, we characterized common scaffolds in the dataset and demonstrated that the chemical space of solutes is extensive and diverse.

## Background & Summary

Solubility is the property of a chemical solute dissolving in a solvent to form a homogeneous system. Solubility depends on the solvent used and the pressure and temperature at which it was measured<sup>1</sup>. While the solubility of solids in the majority of organic solvents plays an essential role in various chemical fields such as organic chemistry, environmental chemistry, organo(photo)catalysis, and petrochemical industry; the solubility data in specified solvents such as water or dimethyl sulfoxide (DMSO) is highly crucial for drug development<sup>2-5</sup>. For example, low solubility in water is attributed with poor bioavailability because compounds that are non-soluble in DMSO could not be used in automatized high-throughput screening (HTS). Generally, it makes them lost for experimental measurements<sup>4</sup>.

There are two general ways to solubility modelling: bottom-up and top-down. Bottom-up approaches are primarily based on molecular simulations<sup>6,7</sup>. Top-down approaches are based on thermodynamic models and quantitative structure property relationship (QSPR) methods<sup>8</sup>. Nowadays, QSPR methods for solubility predictions heavily use machine learning as statistical backbone.<sup>9-11</sup> Machine learning based methods demonstrate high-performance, but rely entirely on the quality and amount of training data<sup>12</sup>. Providing high-quality data is of great importance for further development of top-down approaches.

There are many possible modeling endpoints related to solubility. Solubility can be measured in various solvents and at different temperatures. The interest of pharmaceutical researchers is biased towards the solubility in water and DMSO under normal conditions. However, for organic chemistry, many other solvents can be at the same level of importance. During synthesis and for the purification and crystallization of compounds, one needs to know solubility of compounds in various organic solvents. Our observations indicate that there is lack of solubility data measured in organic solvents, rather than water and DMSO<sup>13,14</sup>. Because solubility is a temperature-dependent property, the solubility data for a number of temperatures is also important. Several solubility datasets currently exist covering a wide range of temperatures and solvents, but the number of measured compounds is not enough for efficient training of machine learning models.<sup>15-17</sup>

The aim of this work is to create a large dataset of organic compounds solubility values to supplement available data for further data-driven model development. For this purpose we searched for and collected data from 771 unique articles from 3 previously selected journals (Figure 1). As a result, organic compounds solubilities in 138 different solvents were collected. This dataset includes solubility values for 839 unique organic compounds under temperatures from 243,15 K to 403,15 K at atmospheric pressure. The complete dataset contains 54273 solubility data records. We named this dataset BigSolDB.

39 We present an open-access dataset of solubility in organic solvents and water in a wide range of temperatures as the main  
40 outcome of our research. We expect it to serve a broad community as a reference dataset for the bench-marking of new  
41 experimental modeling results and also as the data resource for future machine-learning investigations in the chemoinformatics  
42 field, for example, to improve the prediction capability of future models.

43 Humans prefer the visual examination of large volumes of data. Recently, we proposed the parametric t-SNE method as a  
44 universal tool for visualizing chemical data<sup>18,19</sup>. In this research, we provide an interactive web-based tool to visualize the  
45 chemical space of solutes. Using our tool, chemists can explore our dataset and see the solubility charts for each compound  
46 individually. Additionally, we provided the description of significant scaffolds in BigSolDB dataset.

## 47 **Methods**

### 48 **Data collection**

49 The dataset construction was conducted in four steps: (1) choosing appropriate journals, (2) getting data from Crossref API and  
50 filtering papers, (3) extracting data from papers, (4) cleaning and preprocessing data. Our steps are outlined in Figure 1.

### 51 **Journals Selection**

52 In the first step, we selected three journals containing solubility data: *Journal of Chemical & Engineering Data* (ISSN:  
53 1520-5134), *The Journal of Chemical Thermodynamics* (ISSN: 0021-9614), and *Fluid Phase Equilibria* (ISSN: 0378-3812).  
54 These journals focused on publishing compounds' physicochemical data. The motivation behind choosing these journals was  
55 that they systematically publish compounds' solubility data for several solvents at various temperatures.

### 56 **Getting data from Crossref API and filtering papers**

57 In the second step, we searched for the International Standard Serial Number (ISSN) for each journal and, using the Crossref  
58 API (<https://api.crossref.org/>) and obtained a list of Digital Objects Identifiers (DOIs) and titles for each paper.

- 59 • *Journal of Chemical & Engineering Data*: (ISSN: 1520-5134)
- 60 • *The Journal of Chemical Thermodynamics*: (ISSN: 0021-9614)
- 61 • *Fluid Phase Equilibria*: (ISSN: 0378-3812)

62 Using the Crossref bibliographic data, we searched for "solub" subword in titles to pick only papers with solubility data.  
63 We performed manual processing to remove papers that contain "solub" in their titles but do not have the solubility data we  
64 needed. At the final stage, we had 771 papers ready to further processing.

### 65 **Extracting data from papers**

66 We recognized and converted solubility data into a machine-readable format manually. We focused only on mole fraction  
67 solubility under atmospheric pressure. We used OPSIN open-source cheminformatics software ([https://opsin.ch.cam.  
68 ac.uk/](https://opsin.ch.cam.ac.uk/)) to convert IUPAC names of compounds into SMILES representation.

### 69 **Cleaning and preprocessing data**

70 In the chemical literature, solvents can be represented in various ways, which poses a significant challenge for dataset  
71 preparation. To address this issue, we standardized the solvents to eliminate ambiguity. For instance, we converted terms like  
72 "ethyl alcohol" and "methylcarbinol" to "ethanol" while "2-butanol", "butan-2-ol", and "sec-butyl alcohol" were transformed  
73 into "sec-butanol". Subsequently, we employed the OPSIN system to convert the standardized solvent names. Lastly, we  
74 utilized the open-source cheminformatics software RDKit (<https://www.rdkit.org>) to perform SMILES canonization  
75 for each solute compound.

### 76 **Parametric t-SNE visualization**

77 To make exploring our database easier for the researchers, we prepared an interactive web-based tool to visualize the solutes'  
78 chemical space (Figure 2). Our tool helps navigate the 2D scatter plot leveraging several standard operations (zoom areas,  
79 pick compounds, etc.). Also, one can filter compounds by a solvent of interest. The projection is based on the parametric  
80 t-SNE method. The parametric t-SNE model is an artificial neural network that projects molecular structures to 2D (X and  
81 Y) coordinates, keeping the local similarity between molecules. We used the pretrained model from our previous work<sup>18</sup>.  
82 This model was trained on complete ChEMBL data and provided good discriminating ability between organic scaffolds. The  
83 advantage of the parametric t-SNE over the vanilla t-SNE is the determinacy of the model. Because a pretrained model is an  
84 artificial neural network with fixed coefficients, it projects new compounds to known-in-advance regions of chemical space.  
85 Utilizing this feature, one can compare several datasets together, using one of them as the reference dataset.

## 86 Data Records

87 In this research we created BigSolDB: a new broad solubility dataset of a wide range of organic compounds. We provide it as  
88 a comma-separated values (csv) file. BigSolDB contains 54273 individual solubility values, 830 unique molecules and 138  
89 individual solvents for the temperature range from 243.15 to 403.15K at atmospheric pressure. The dataset is freely accessible  
90 in Zenodo data sharing platform<sup>20</sup>.

91 The file has five columns:

- 92 • SMILES — SMILES representation of a dissolved compound
- 93 • T,K — Temperature in Kelvin
- 94 • Solubility — an experimental solubility value (mole fraction)
- 95 • Solvent — a name of the solvent
- 96 • SMILES\_Solvent — SMILES representation of the solvent
- 97 • Source — a data source for the record

98 Our interactive visualisation tool is freely accessible: <http://bigsoldbapp.cheminfo.space>. The project page  
99 is located: <https://bigsoldb.cheminfo.space>.

## 100 Data Analysis

101 One can see that top-8 solvents (ethanol, methanol, isopropanol, ethyl acetate, acetone, n-propanol, water, n-butanol, acetonitrile)  
102 cover 57% of data measurements in total, wherein water (3168 records, 5.8%) is only at the seventh place. Probably, this reflex  
103 the most common solvents that are using in organic chemistry. The similar picture we have for the individual compounds.  
104 The majority of experiments was performed in alcohols, not in water. However, this bias seems undesirable for the medicinal  
105 chemistry, because, the water solubility for drug candidates is of greater importance.

106 Following The United States Pharmacopeia (USP, a nonprofit organization committed to establishing standards for medicines,  
107 food, dietary supplements, etc.) approach, compounds can be divided into seven categories with respect to their solubilities  
108 ratios<sup>21</sup>. We converted the molar fractions into molar ratios under fixed temperature (293.15 K) and attributed them to one of  
109 seven USP classes (Table 2). The pie chart demonstrating the distribution of USP classes is given in Figure 14. One can see that  
110 six out of seven classes are well-represented in the dataset, the exception is "Very soluble" class which is presented only by 39  
111 records (0.98% of total records at 293.15 K).

112 The data distribution across solvents was presented as a set of barplots, where each bar corresponds to the number of data  
113 rows (Figure 3, Figure 4) or the number of individual compounds (Figure 5, Figure 6) in the dataset. The top 50 of the most  
114 popular solvents were analysed separately (Figure 3, Figure 5).

115 The values of mole fraction solubilities shown as the histogram in Figure 7 are disturbed from  $5.5 * 10^{-11}$  to  $9.506 * 10^{-1}$ .  
116 Figure 8 shows the distribution of the number of solubility observations measured at various temperatures.

117 Top-9 compounds that have the largest number of records are depicted in Figure 12. One can see, that 8 out of 9 of these  
118 compounds contain several functional groups with an oxygen atom (caboxyl-, nitro-, ester-, etc.).

119 Temperature coverage of the solubilities measurements in the 30 most popular solvents was presented in Figure 9. Right  
120 heat bar shows the counts of data rows for solubilities under corresponding temperatures. Temperatures to show were chosen  
121 according to the histogram in Figure 8. This figure shows the results for the temperature range with the largest number of  
122 records (from 278.15 to 323.15 K) in the dataset.

123 Also, data coverage was represented as a heatmap in Figure 10. We chose SMILES to show as the top 50 of the mostly  
124 described compounds in the dataset. And we analysed their solubilities in the 30 of the most popular solvents. Right heat bar  
125 represents to the data rows number of the solubilities of described compounds in shown solvents. As shown in the Figure 11,  
126 the majority of presented compounds occurs in the dataset only once. Because of this fact and the lack of the experimental  
127 errors information in the original articles, no statistical solubility data validation would be performed on the dataset.

128 We also compared BigSolDB with existing ones<sup>15,17</sup> and we present the results on Table 1. BigSolDB contains ten times  
129 more solubility values than previous datasets and three times more unique molecules in similar temperature range.

## 130 Chemical space of BigSolDB

131 To reveal the chemical space of BigSolDB, we prepared parametric t-SNE projections of BigSolDB and Tox21 dataset<sup>22</sup>  
132 (reference dataset) on the same plot (Figure 13). One should regard that the number of compounds in Tox21 dataset is larger than  
133 in BigSolDB. So, we had have to set a lower opacity value for Tox21 data points to prevent the visual suppression of BigSolDB

134 points. One can see that the overall distributions of chemical compounds are similar; both Tox21 and BigSolDB are broad and  
135 diverse. We highlighted and described several representative clusters on this map. The cluster (A) comprises "simple" aliphatic  
136 compounds. The cluster (C) is represented by aromatic compounds with the dominance of substituted benzenes. cluster  
137 (H) comprises of aliphatic heterocycles. The groups of steroids (F) and sugar derivatives (G) are well-recognisable. There  
138 are two clusters (D) and (E) which consists of two established clusters: N-phenyl-1-thia-3-azaspiro[5.5]undec-2-en-2-amine  
139 derivatives and Desethylenenorfloxacin derivatives. However, BigSolDB is notably enriched with sulfonamide derivatives (B)  
140 been compared with Tox21 dataset.

## 141 Usage Notes

142 The given dataset can be used both to view individual solubility values and to train machine learning models. To create  
143 fingerprints for machine learning models SMILES representations of dissolved compounds could be converted to an RDKit  
144 molecular object. The presented data is available under CC-BY 4.0 licence.

## 145 Code Availability

146 The Jupyter Notebook scripts for downloading data from Crossref API and plotting the figures can be found in [https:](https://github.com/Cemich/BigSolDB)  
147 [//github.com/Cemich/BigSolDB](https://github.com/Cemich/BigSolDB).

## 148 References

- 149 1. Boobier, S., Osbourn, A. & Mitchell, J. B. O. Can human experts predict solubility better than computers? *J. Cheminform-*  
150 *matics* **9**, 63, [10.1186/s13321-017-0250-y](https://doi.org/10.1186/s13321-017-0250-y) (2017).
- 151 2. Delaney, J. S. Esol: estimating aqueous solubility directly from molecular structure. *J. Chem. Inf. Comput. Sci.* **44**,  
152 1000–1005, [10.1021/ci034243x](https://doi.org/10.1021/ci034243x) (2004).
- 153 3. Sorkun, M. C., Khetan, A. & Er, S. Aqsolddb, a curated reference set of aqueous solubility and 2d descriptors for a diverse  
154 set of compounds. *Sci. Data* **6**, 143, [10.1038/s41597-019-0151-1](https://doi.org/10.1038/s41597-019-0151-1) (2019).
- 155 4. Tetko, I. V. *et al.* Development of dimethyl sulfoxide solubility models using 163 000 molecules: Using a domain  
156 applicability metric to select more reliable predictions. *J. Chem. Inf. Model.* **53**, 1990–2000, [10.1021/ci400213d](https://doi.org/10.1021/ci400213d) (2013).
- 157 5. Baybekov, S. *et al.* Dmso solubility assessment for fragment-based screening. *Molecules* **26**, [10.3390/molecules26133950](https://doi.org/10.3390/molecules26133950)  
158 (2021).
- 159 6. Hossain, S., Kabedev, A., Parrow, A., Bergström, C. A. & Larsson, P. Molecular simulation as a computational pharma-  
160 ceutics tool to predict drug solubility, solubilization processes and partitioning. *Eur. J. Pharm. Biopharm.* **137**, 46–55,  
161 <https://doi.org/10.1016/j.ejpb.2019.02.007> (2019).
- 162 7. Bjelobrk, Z., Mendels, D., Karmakar, T., Parrinello, M. & Mazzotti, M. Solubility prediction of organic molecules  
163 with molecular dynamics simulations. *Cryst. Growth & Des.* **21**, 5198–5205, [10.1021/acs.cgd.1c00546](https://doi.org/10.1021/acs.cgd.1c00546) (2021). [https:](https://doi.org/10.1021/acs.cgd.1c00546)  
164 [//doi.org/10.1021/acs.cgd.1c00546](https://doi.org/10.1021/acs.cgd.1c00546).
- 165 8. Bouillot, B., Teychené, S. & Biscans, B. An evaluation of thermodynamic models for the prediction of drug and drug-like  
166 molecule solubility in organic solvents. *Fluid Phase Equilibria* **309**, 36–52, <https://doi.org/10.1016/j.fluid.2011.06.032>  
167 (2011).
- 168 9. Kuentz, M. & Bergström, C. A. Synergistic computational modeling approaches as team players in the game of solubility  
169 predictions. *J. Pharm. Sci.* **110**, 22–34, <https://doi.org/10.1016/j.xphs.2020.10.068> (2021).
- 170 10. Mullins, E., Liu, Y. A., Ghaderi, A. & Fast, S. D. Sigma profile database for predicting solid solubility in pure and mixed  
171 solvent mixtures for organic pharmacological compounds with cosmo-based thermodynamic methods. *Ind. & Eng. Chem.*  
172 *Res.* **47**, 1707–1725, [10.1021/ie0711022](https://doi.org/10.1021/ie0711022) (2008).
- 173 11. Ran, Y., Jain, N. & Yalkowsky, S. H. Prediction of aqueous solubility of organic compounds by the general solubility  
174 equation (gse). *J. Chem. Inf. Comput. Sci.* **41**, 1208–1217, [10.1021/ci010287z](https://doi.org/10.1021/ci010287z) (2001).
- 175 12. Artrith, N. *et al.* Best practices in machine learning for chemistry. *Nat. Chem.* **13**, 505–508, [10.1038/s41557-021-00716-z](https://doi.org/10.1038/s41557-021-00716-z)  
176 (2021).
- 177 13. Abraham, M. H. *et al.* Prediction of solubility of drugs and other compounds in organic solvents. *J. Pharm. Sci.* **99**,  
178 1500–1515, <https://doi.org/10.1002/jps.21922> (2010).
- 179 14. Lim, H. & Jung, Y. Delfos: deep learning model for prediction of solvation free energies in generic organic solvents. *Chem.*  
180 *Sci.* **10**, 8306–8315, [10.1039/C9SC02452B](https://doi.org/10.1039/C9SC02452B) (2019).

- 181 **15.** Ye, Z. & Ouyang, D. Prediction of small-molecule compound solubility in organic solvents by machine learning algorithms.  
182 *J. Cheminformatics* **13**, 98, [10.1186/s13321-021-00575-3](https://doi.org/10.1186/s13321-021-00575-3) (2021).
- 183 **16.** Boobier, S., Hose, D. R. J., Blacker, A. J. & Nguyen, B. N. Machine learning with physicochemical relationships: solubility  
184 prediction in organic solvents and water. *Nat. Commun.* **11**, 5753, [10.1038/s41467-020-19594-z](https://doi.org/10.1038/s41467-020-19594-z) (2020).
- 185 **17.** Vermeire, F. H., Chung, Y. & Green, W. H. Predicting solubility limits of organic solutes for a wide range of solvents and  
186 temperatures. *J. Am. Chem. Soc.* **144**, 10785–10797, [10.1021/jacs.2c01768](https://doi.org/10.1021/jacs.2c01768) (2022).
- 187 **18.** Karlov, D. S., Sosnin, S., Tetko, I. V. & Fedorov, M. V. Chemical space exploration guided by deep neural networks. *RSC*  
188 *Adv.* **9**, 5151–5157, [10.1039/C8RA10182E](https://doi.org/10.1039/C8RA10182E) (2019).
- 189 **19.** Andronov, M., Fedorov, M. V. & Sosnin, S. Exploring chemical reaction space with reaction difference fingerprints and  
190 parametric t-sne. *ACS Omega* **6**, 30743–30751, [10.1021/acsomega.1c04778](https://doi.org/10.1021/acsomega.1c04778) (2021).
- 191 **20.** Krasnov, L., Mikhaylov, S., Fedorov, M. & Sosnin, S. Bigsolddb: solubility dataset of compounds in organic solvents and  
192 water in a wide range of temperatures. *zenodo* <https://doi.org/10.5281/zenodo.6809668> (2022).
- 193 **21.** Convention, U. S. P. *USP35 NF30, 2012: U. S. Pharmacopoeia National Formulary*. The United States pharmacopeia  
194 (United States Pharmacopeial, 2011).
- 195 **22.** Huang, R. *et al.* Tox21challenge to build predictive models of nuclear receptor and stress response pathways as mediated  
196 by exposure to environmental chemicals and drugs. *Front. Environ. Sci.* **3**, [10.3389/fenvs.2015.00085](https://doi.org/10.3389/fenvs.2015.00085) (2016).

## 197 Acknowledgements

198 We used Figma (<https://www.figma.com/>) and Seaborn (<https://seaborn.pydata.org/>) to create all draw-  
199 ings in this paper. We used Dash (<https://dash.plotly.com/>) and plotly (<https://plotly.com/>) to create the  
200 visualisation tool.

## 201 Author contributions statement

202 L.K. and S.M. collected the dataset and analysed the results. S.S. implemented the tool for the visualisation of the dataset and  
203 provided the description of the chemical space. L.K., S.M. and S.S. wrote the draft of the manuscript, and M.V.F. revised it.  
204 S.S. supervised the project. All authors reviewed the manuscript.

## 205 Competing interests

206 The authors declare no competing interests.

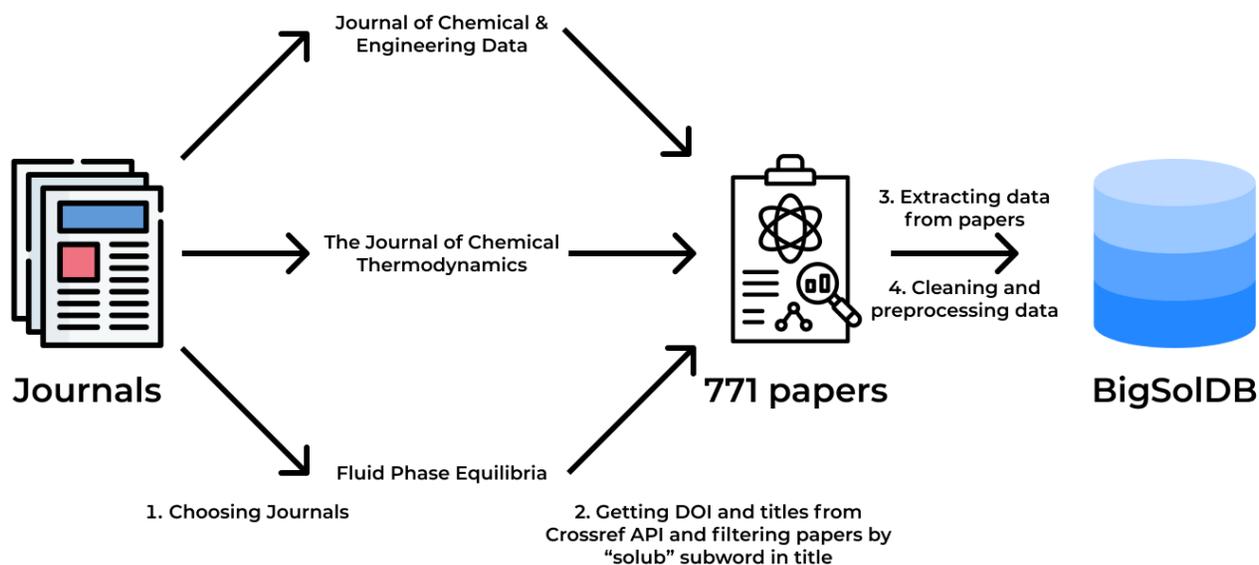
## 207 Figures & Tables

Dataset	Number of records	SMILES	Solvents	Temperature range	Sources
Dataset from <sup>17</sup>	4953	116	98	243.15-363.7	105
Dataset from <sup>15</sup>	5081	266	123	213.15-413.15	—
BigSolDB	54273	839	138	243.15-403.15	771

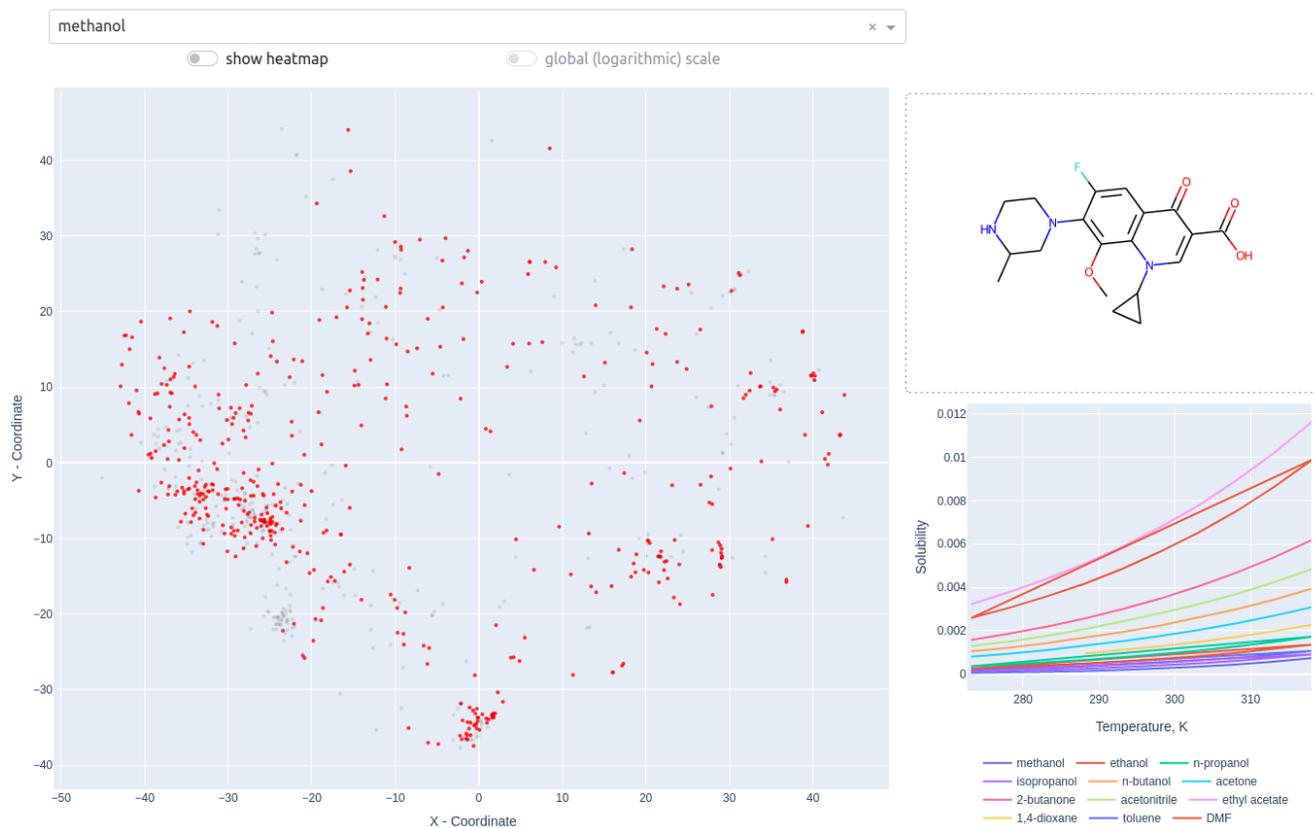
**Table 1.** Comparison of BigSolDB dataset with the previous ones.

Range of molar ratio	Class
0-1	Very soluble
1-10	Freely soluble
10-30	Soluble
30-100	Sparingly soluble
100-1000	Slightly soluble
1000-10000	Very slightly soluble
10000-∞	Practically insoluble

**Table 2.** USP solubility classes.



**Figure 1.** The process of creating a dataset.



**Figure 2.** The demonstration of our solute chemical space visualization tool. On this figure, the particular solvent is selected "methanol", so the only compounds that have methanol solubility data are highlighted red. The rest of structures are faded grey.

Top 50 of the most popular solvents

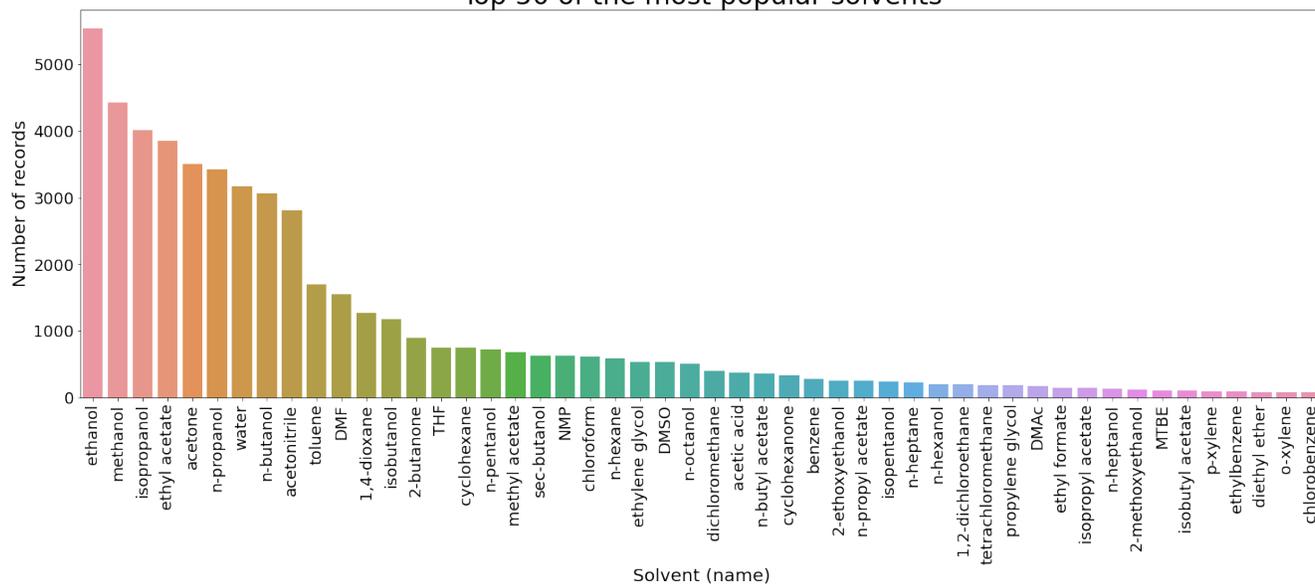


Figure 3. Barplot of the data rows number for the 50 most popular solvents.

Data rows counts in other solvents

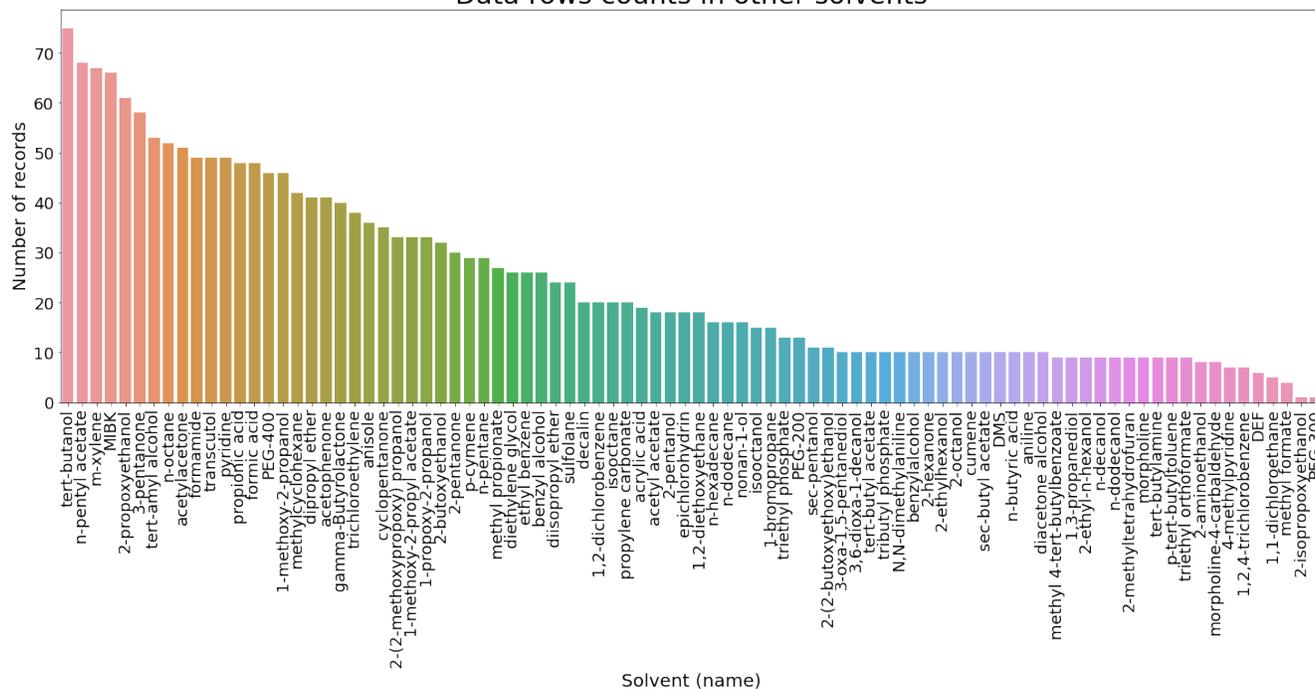


Figure 4. Barplot of the data rows number for the other solvents.

Compounds countplot through the most popular solvents

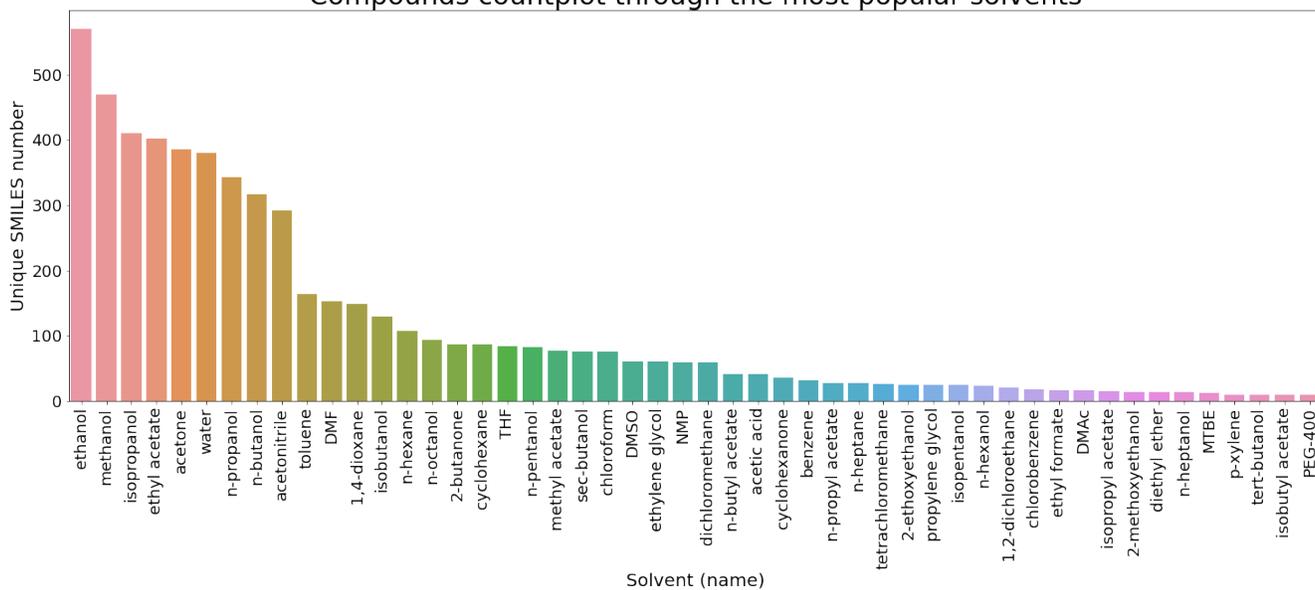


Figure 5. The number of compounds which were dissolved in the 50 most popular solvents.

Compounds countplot through the most popular solvents

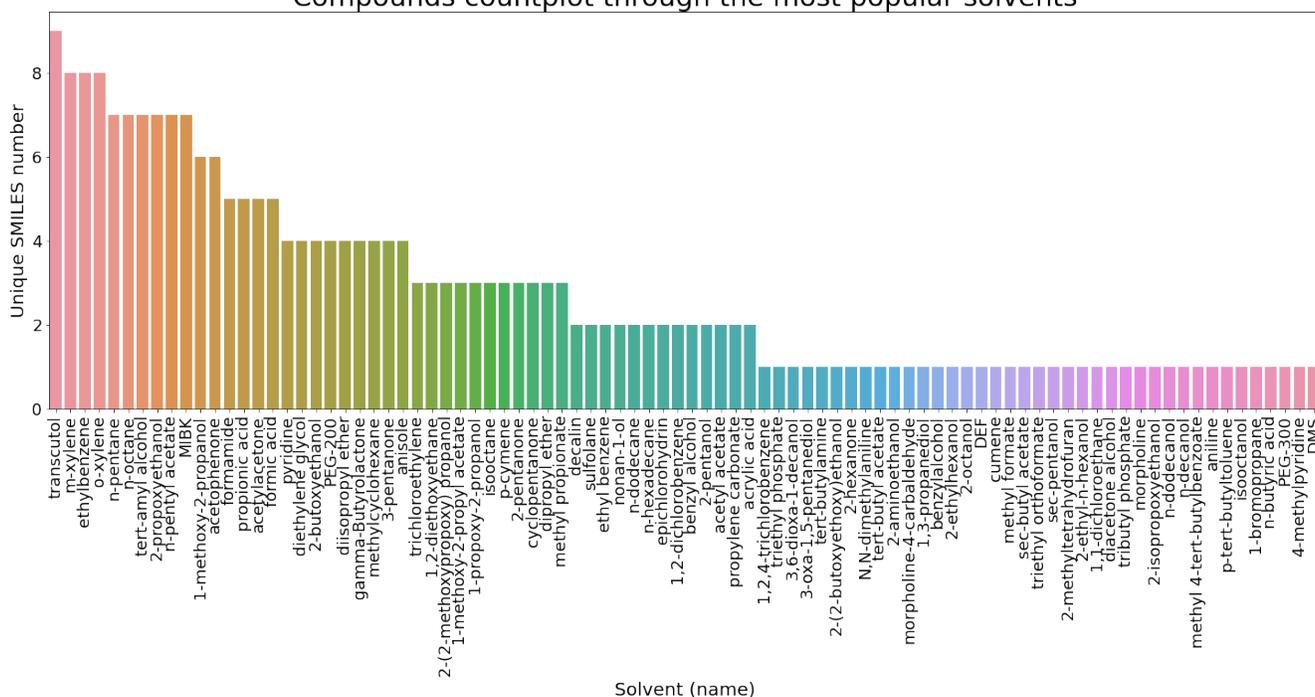
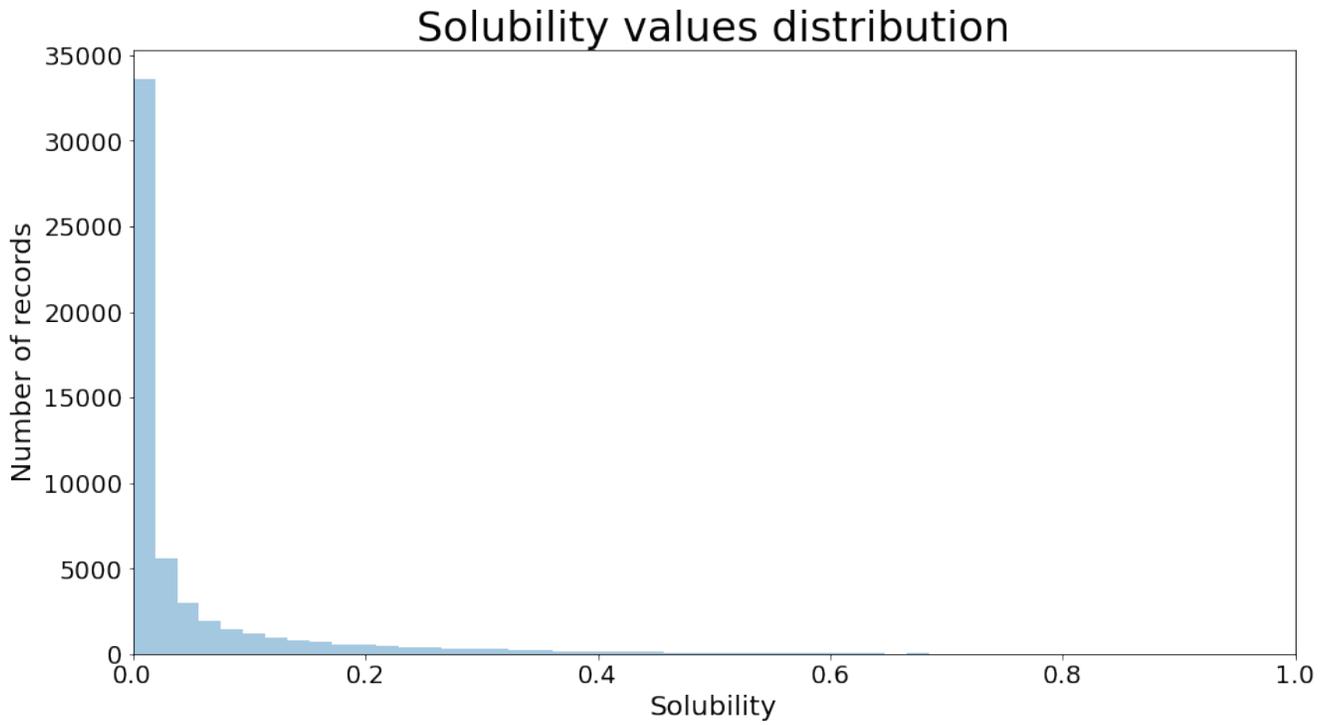
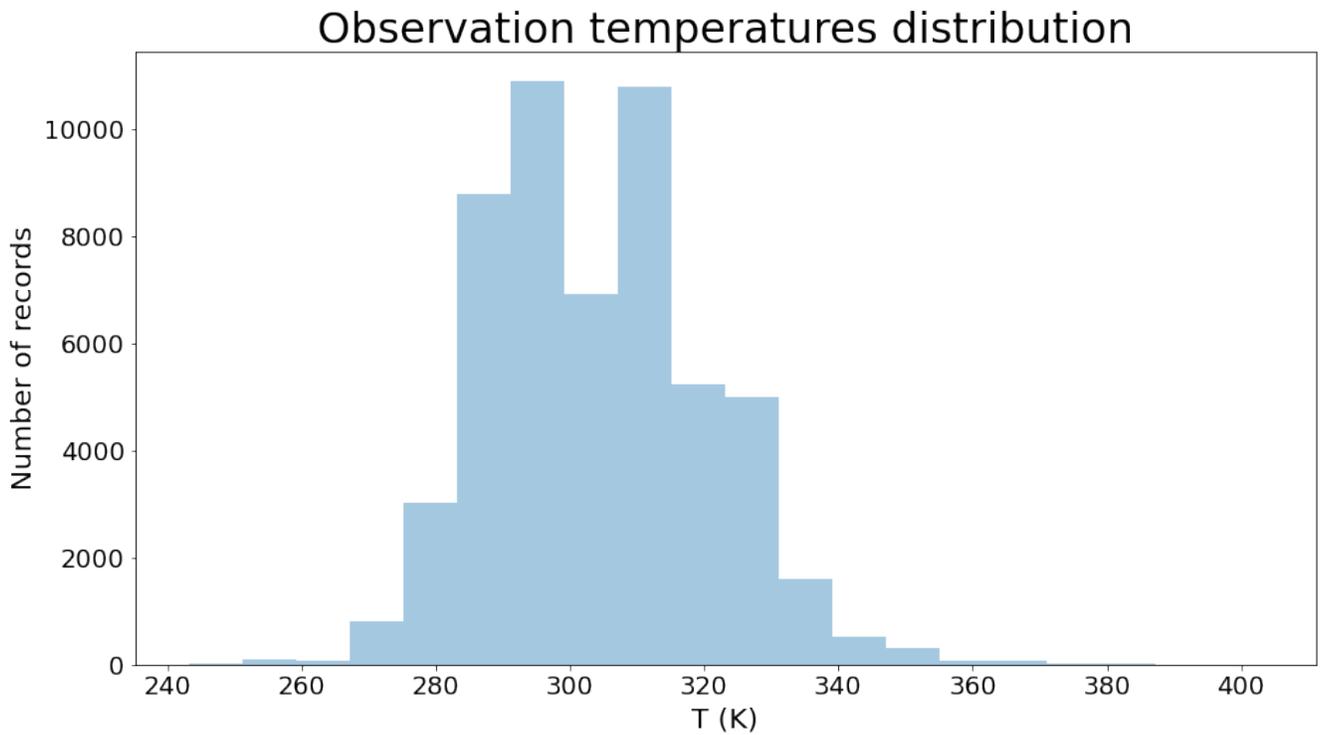


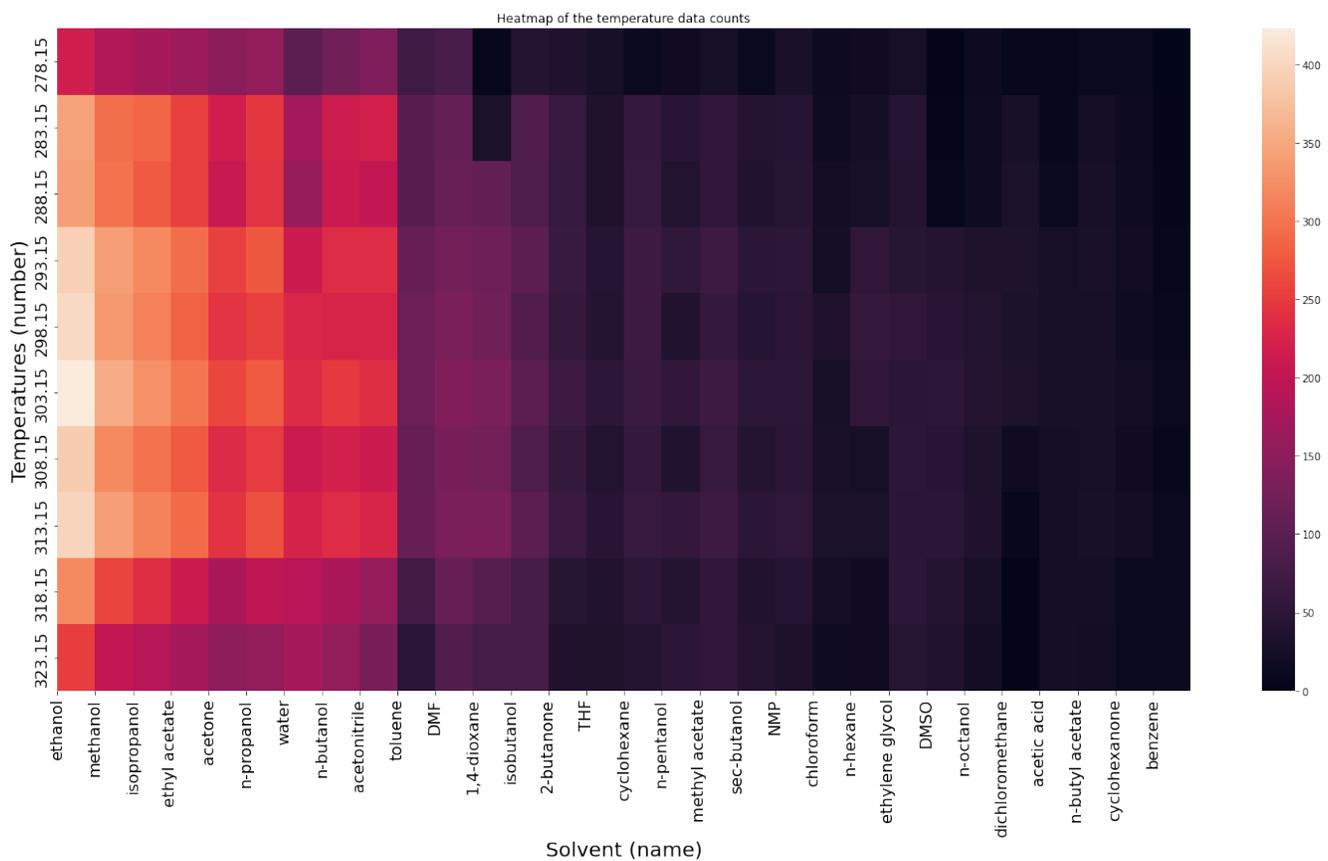
Figure 6. The numbers of compounds which were dissolved in the other solvents.



**Figure 7.** Histogram of all solubility values in the dataset.



**Figure 8.** Histogram of all available temperatures of solubility observations.



**Figure 9.** Heatmap of the number of solubilities measured in the 30 most popular solvents at the most popular temperatures.

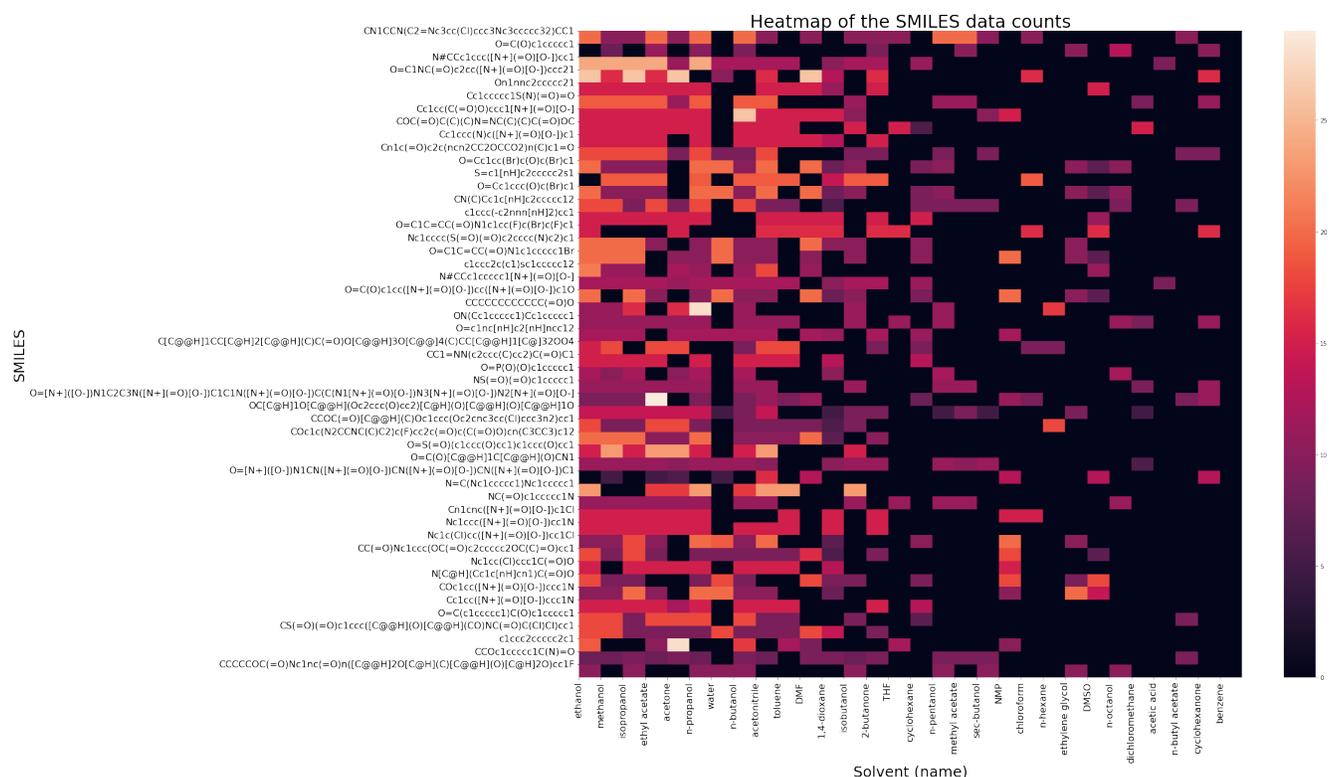


Figure 10. Heatmap of the number of solubilities measured in the 30 most popular solvents for the compounds with the largest number of records.

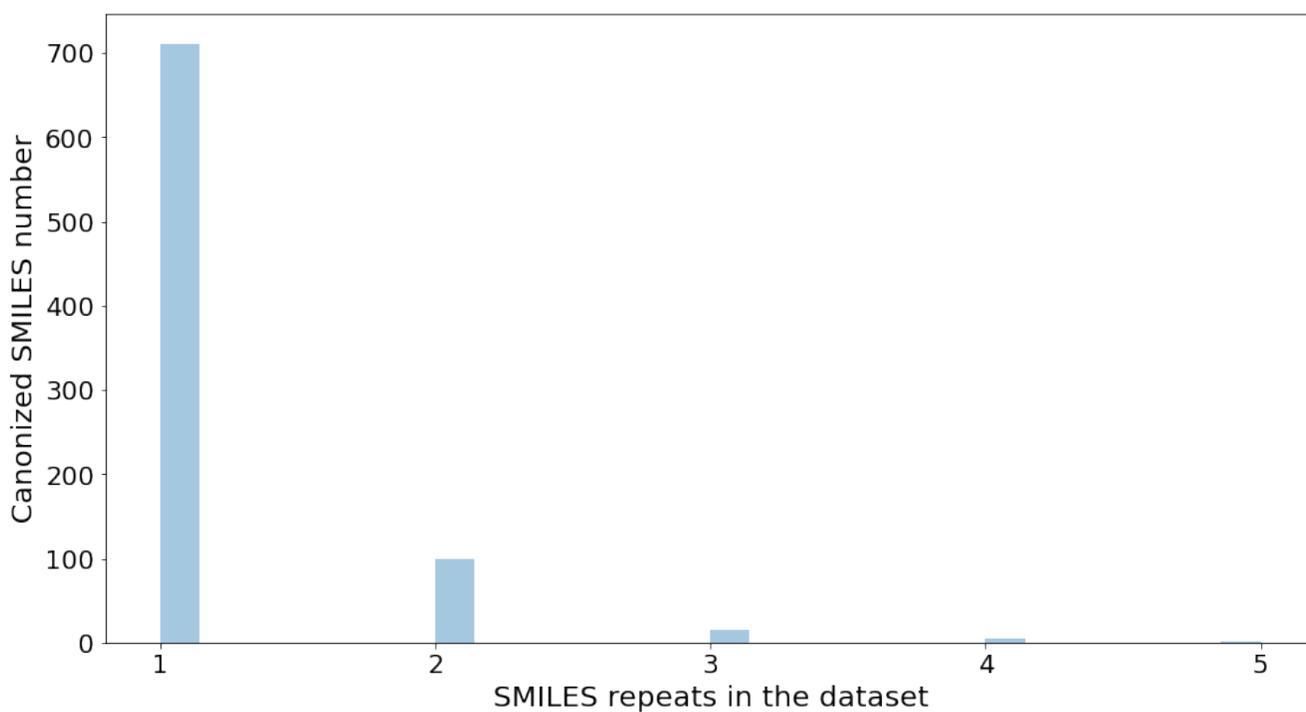
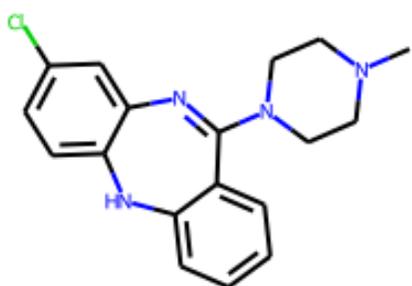
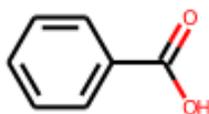


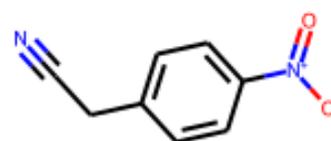
Figure 11. Countplot of the SMILES occurrence in the dataset.



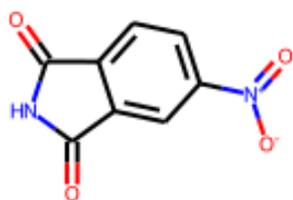
N of records: 300



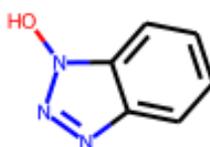
N of records: 278



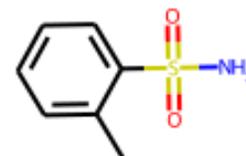
N of records: 245



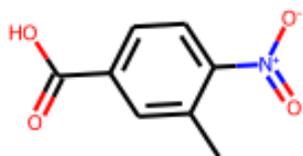
N of records: 239



N of records: 236



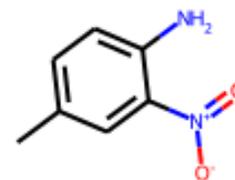
N of records: 230



N of records: 224

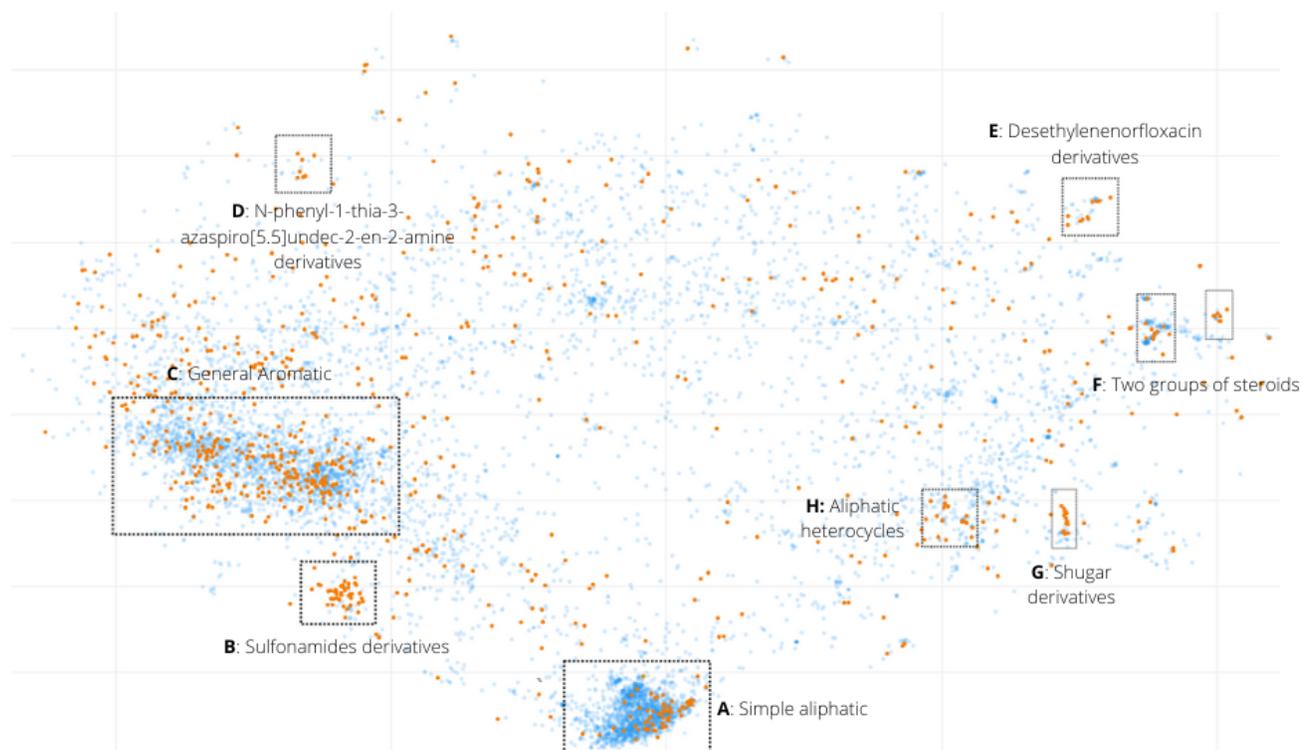


N of records: 215

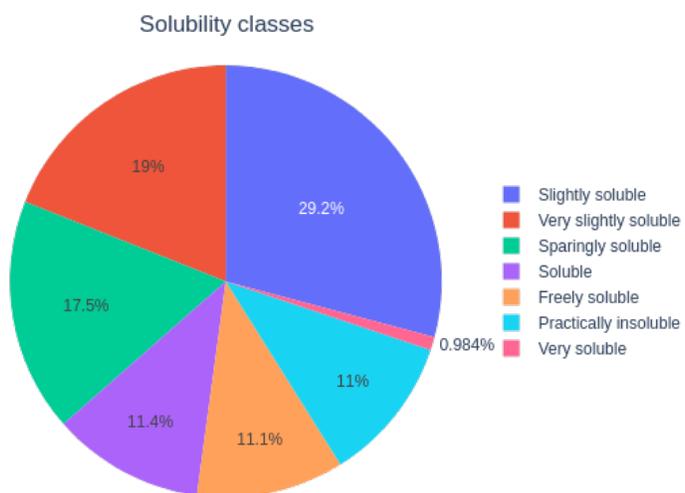


N of records: 214

**Figure 12.** Top-9 of molecules sorted by the maximal number of records in the dataset.



**Figure 13.** The comparison between Tox21 dataset and BigSolDB. Blue dots represent structures from Tox21 dataset, orange ones represent structures from BigSolDB



**Figure 14.** The distribution of USP solubility classes at 293.15 K