# An example to use MSC to analyse the predictions from a machine learning model

Jan-Mathis Hein
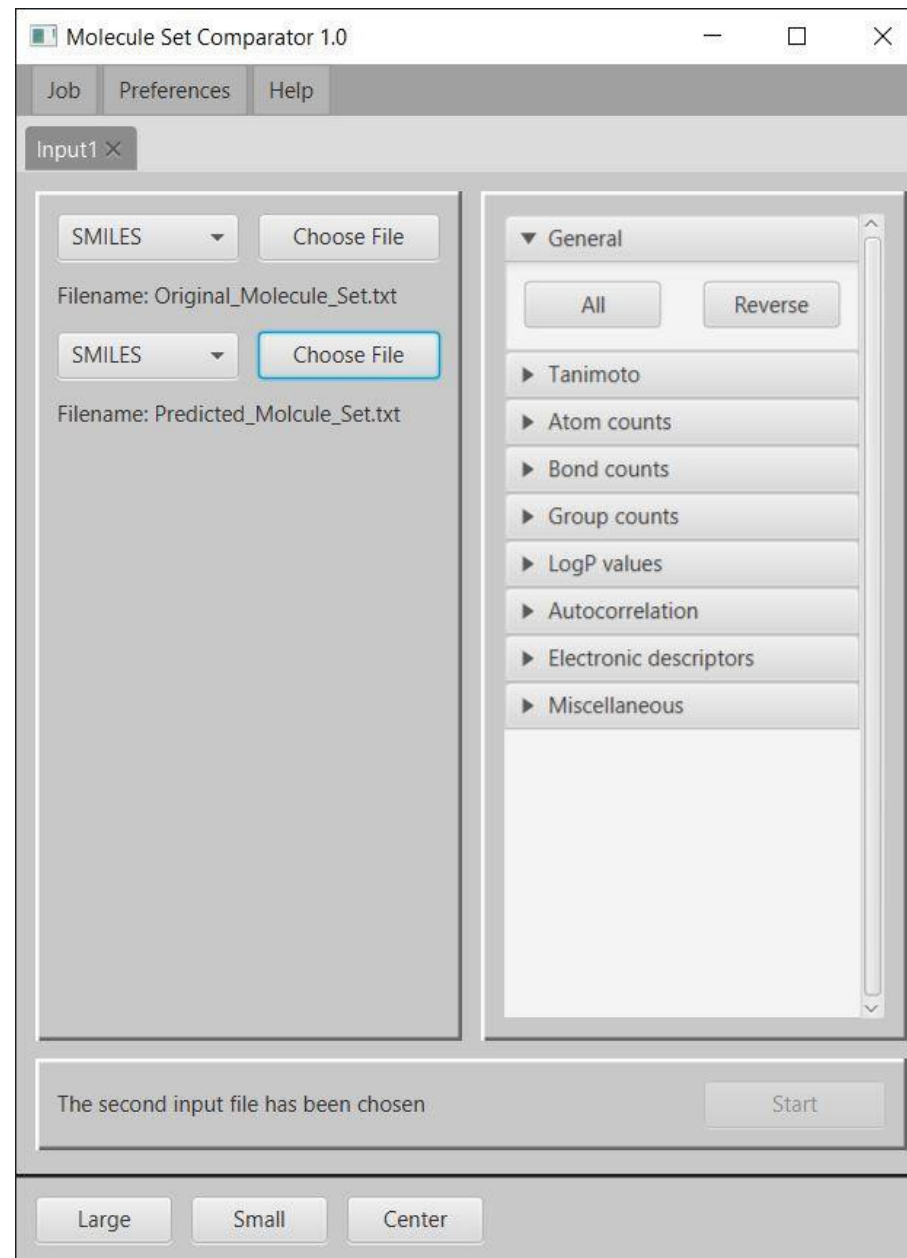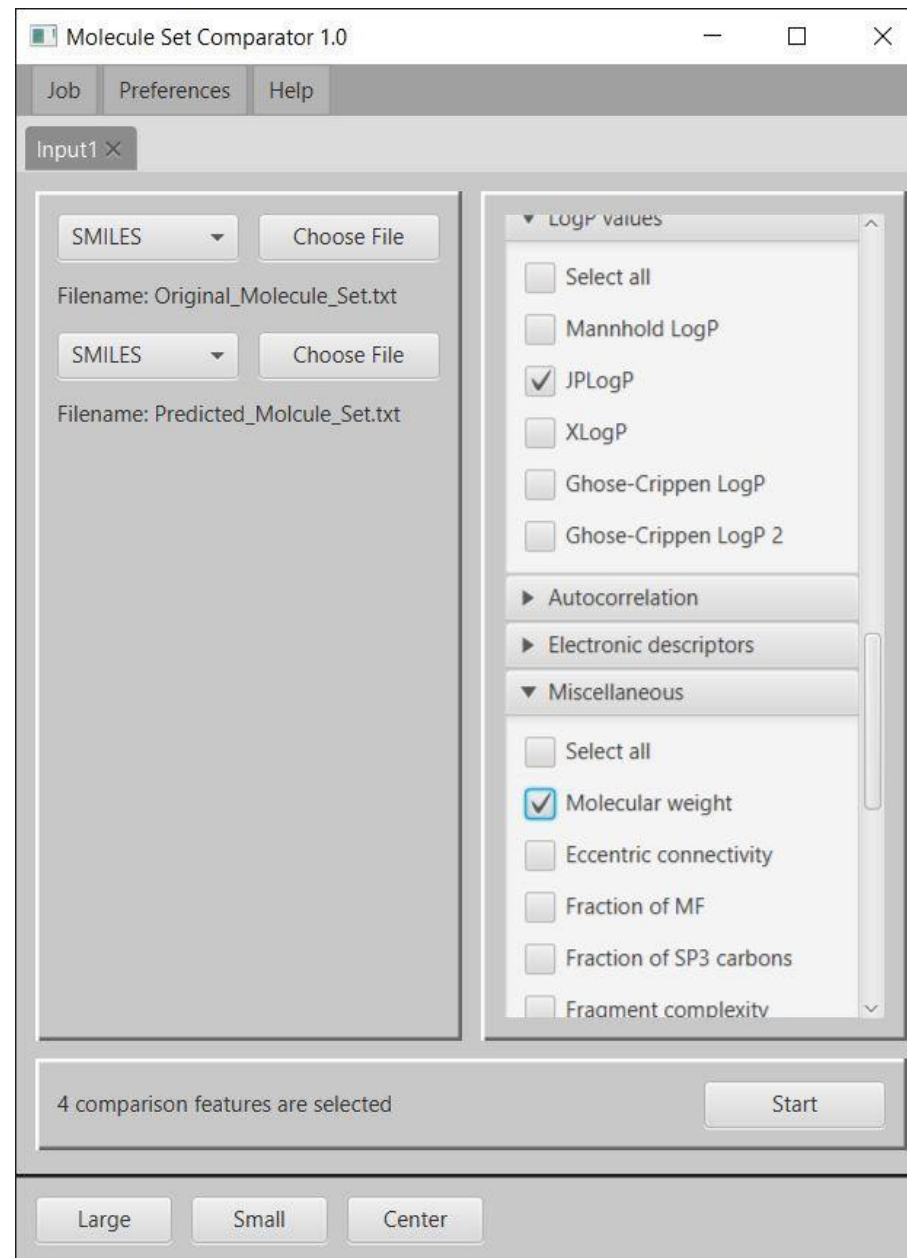
# Initializing the job

- In this tutorial we demonstrate how the MSC can be used to evaluate predictions from a machine learning model

- First, the application is started, and a new job is created.

- Then the original and predicted molecule sets are loaded into the MSC.

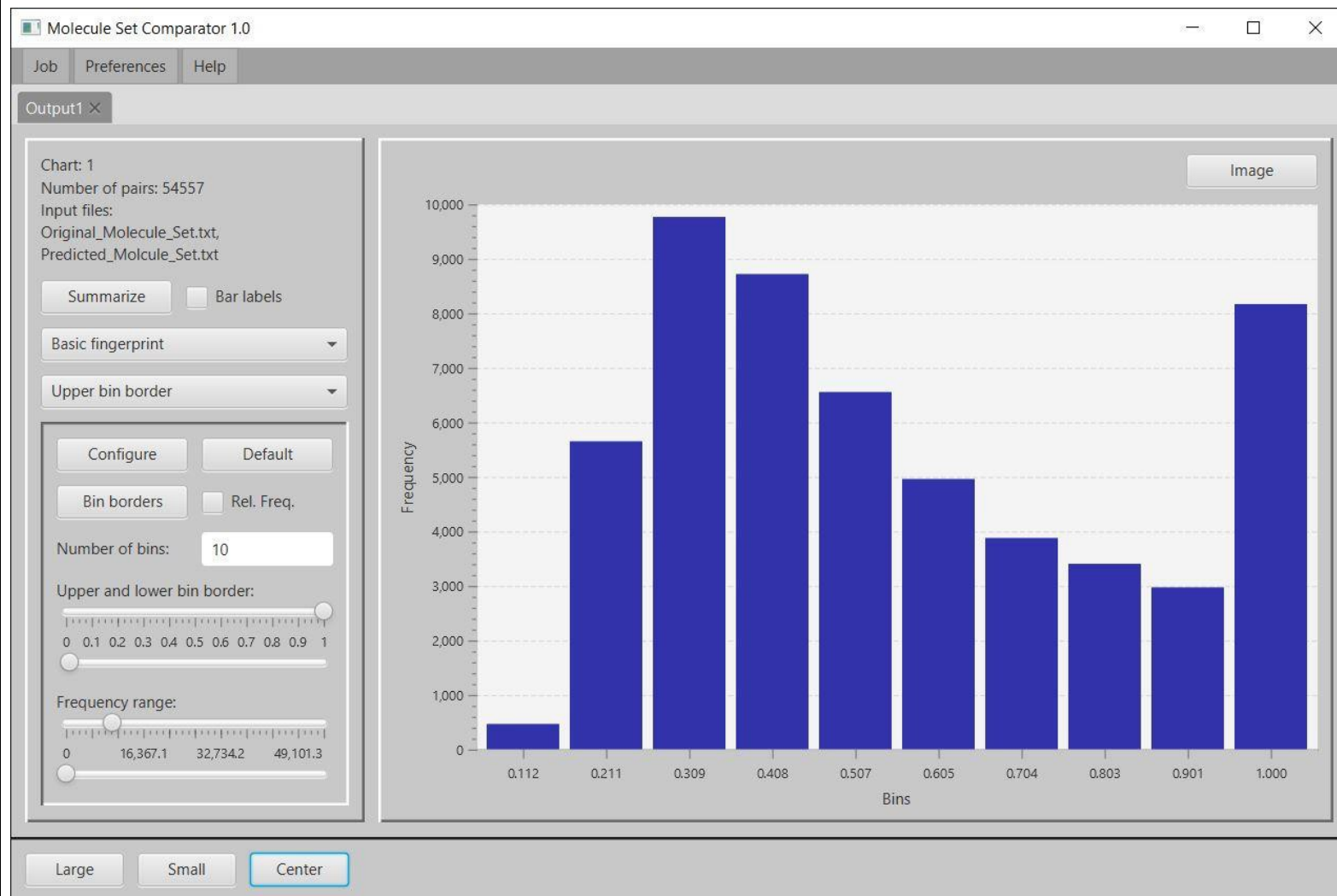- Here, both sets are encoded in the SMILES format.

# Initializing the job

- Next, a set of molecular descriptors can be selected to compare the molecule sets based on the user's need.

- Here we have selected the Tanimoto coefficient with the basic CDK fingerprint, the atom count, the JPLogP value, and the molecular weight.  For the last three descriptors, each pair is compared based on the difference between the two descriptors.

- Once the job gets started, the molecule sets are compared pairwise (the first original molecule to the first predicted molecule and so on)
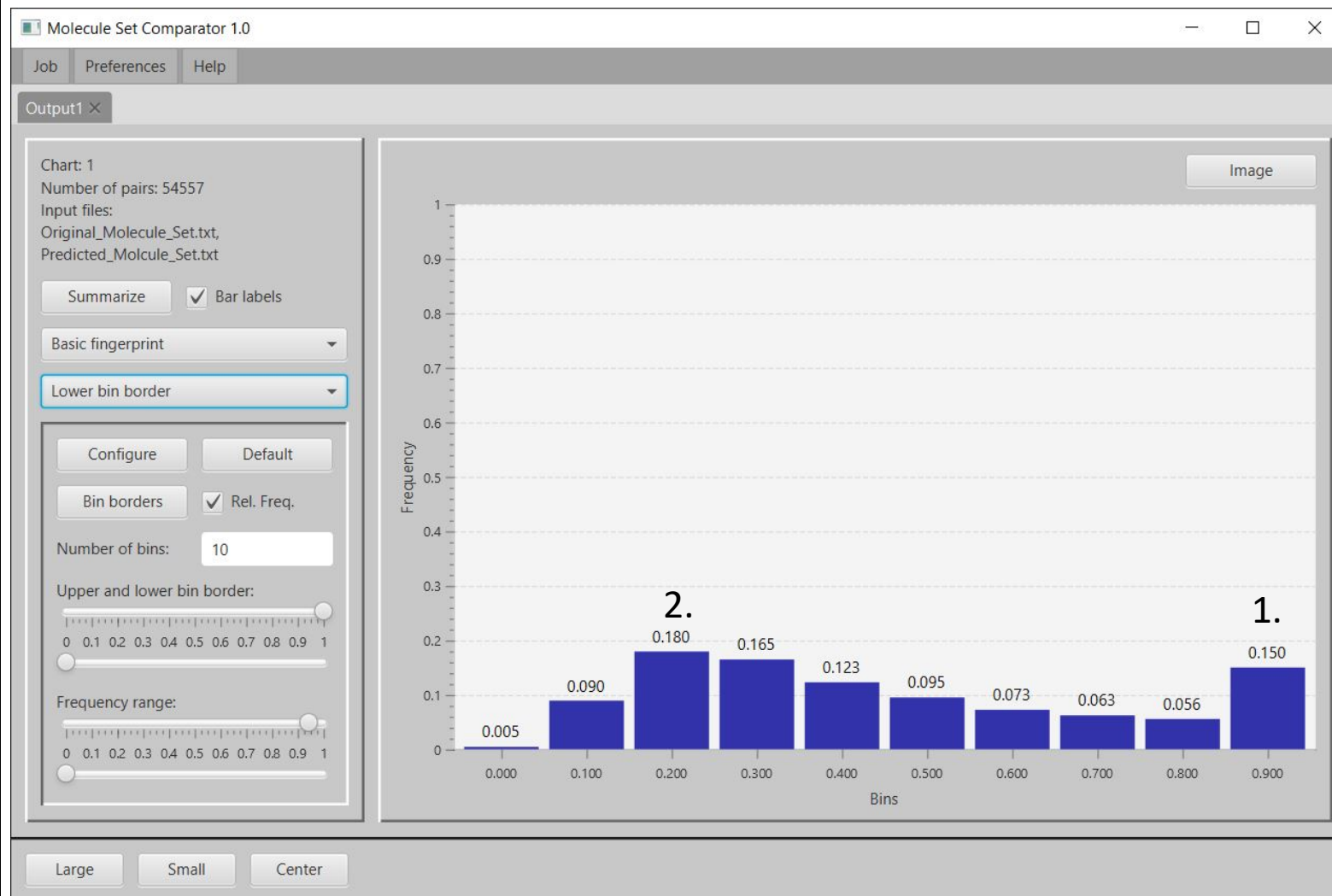
# Analysing the results

- Once the execution completes, a separate tab with the results will pop up.

- The histogram shows the absolute frequencies of molecule pairs for which the Tanimoto coefficient rests in a specific interval. On the x-axis, the intervals are labelled with their upper border (so the last interval goes from 0.901 to 1.0)

- The first step towards analysing these results is to configure the histogram according to your needs.
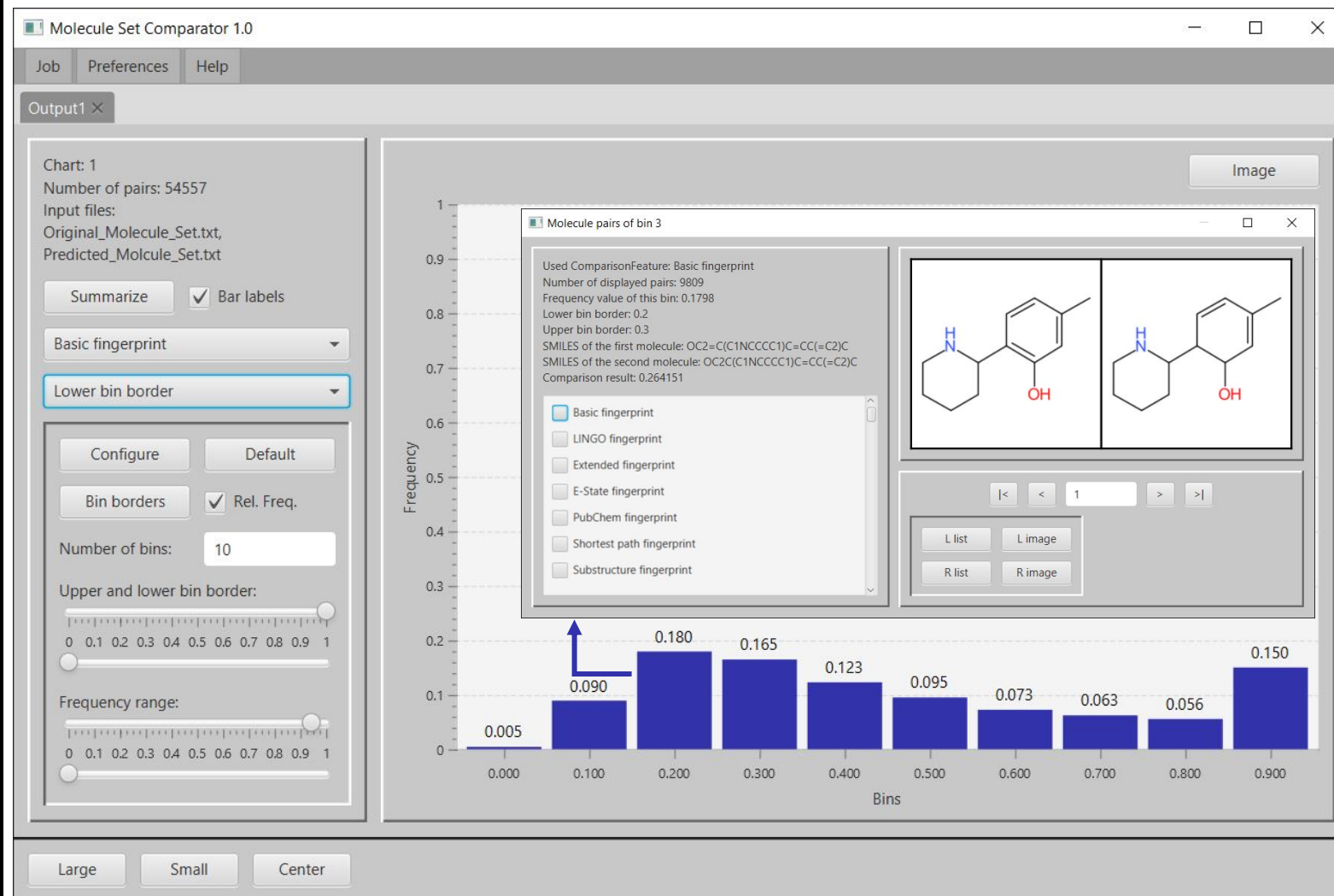
# Analysing the results

- Instead of absolute frequencies, relative ones are now displayed.

- Every bar is labelled with its frequency.

- The y-axis now ranges from 0 to 1

- The x-axis labels now display the lower border of each interval (so the first interval goes from 0.0 to 0.1) and now go from 0.0 to 0.9

- Around 15% of the original set was predicted very well by the algorithm (1.). But many molecules were predicted poorly (2.)

# The detail window

- The poorly predicted molecules can be examined more closely by clicking on one of the bins, for example the third bin (see arrow).

- By clicking on the bin, the detail window for that bin will be opened. Here one can browse through the list of input-output-pairs of the bin.

- It is immediately visible that the first molecule pair doesn't look that different. This intuitive rating contradicts the low Tanimoto coefficient of 0.26
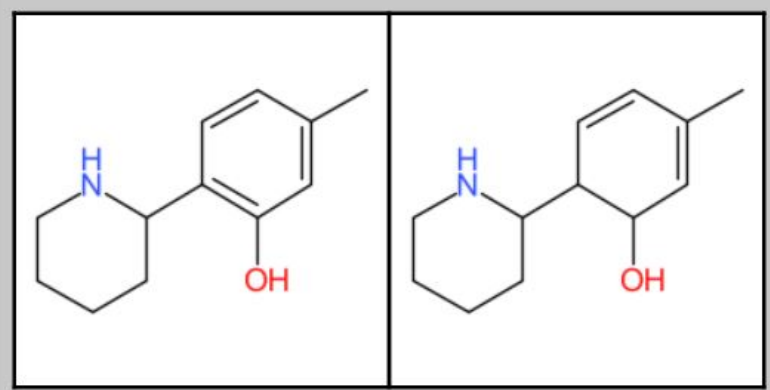
# The detail window

- This perceived contradiction can be solved by calculating additional Tanimoto coefficients with different fingerprints.

- For example, the PubChem fingerprint yields a Tanimoto value of 0.6 for the same molecule pair which seems to be a lot more reasonable.

- This is because the Tanimoto coefficient is very dependent on the fingerprint.

- To get a clearer picture we could run a separate job with a different fingerprint.
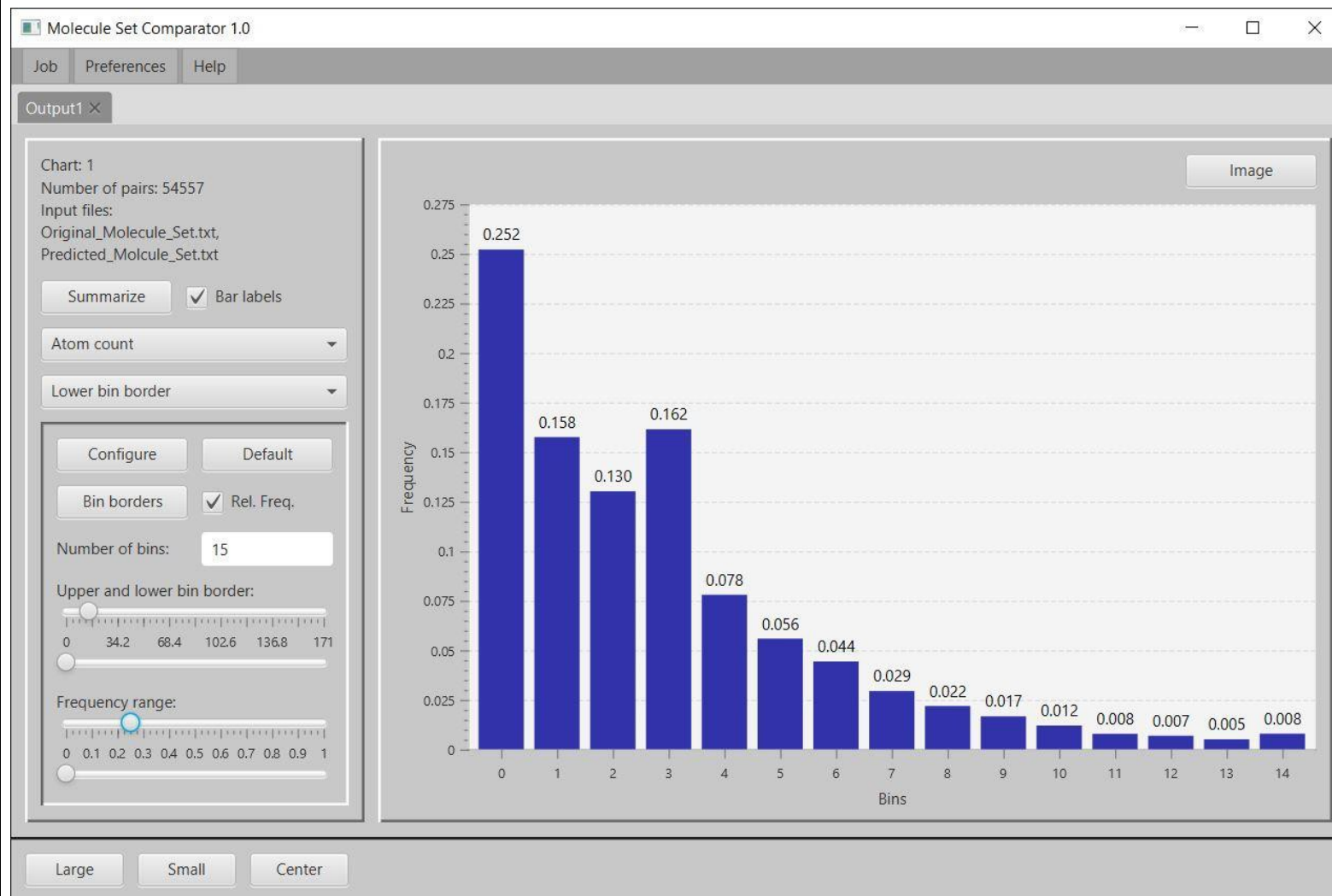
# Other descriptors

- In the output tab, we can quickly switch between the different calculated descriptors.

- Here the histogram of the atom counts descriptor is displayed. The intervals are labelled with their lower border (so the first interval goes from 0 to 17, the second from 18 to 34 and so on)

- If you do not configure the display, the histogram can be uninformative.
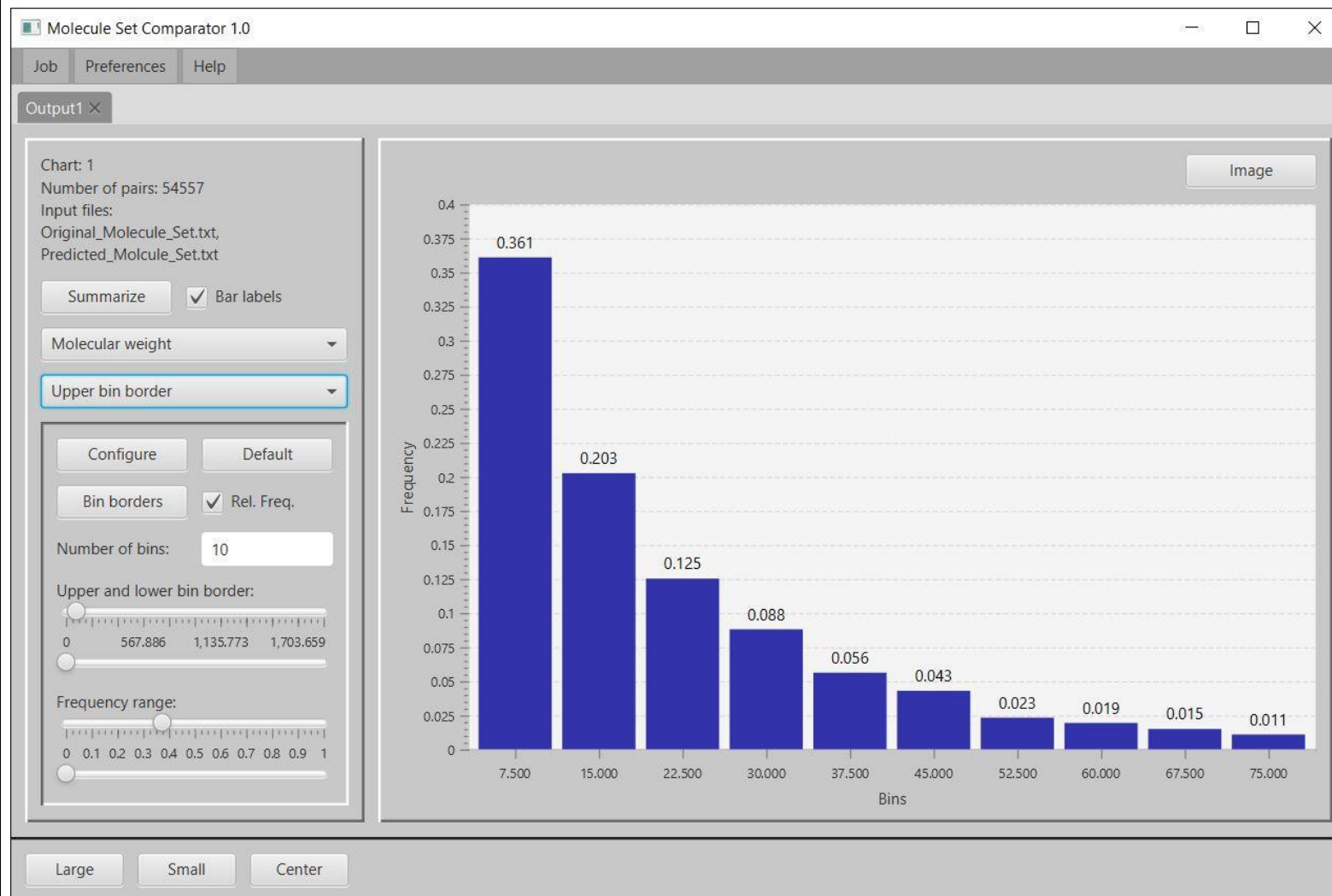
# Other descriptors

- After configuring, the histogram conveys a lot more information.

- Around 25% of the molecules that were predicted have the right number of atoms.

- The histogram shows a "nice" distribution except for the 4th bar.

- This anomaly is likely to be caused by the addition of one "C" atom which also leads to two additional "H" atoms.
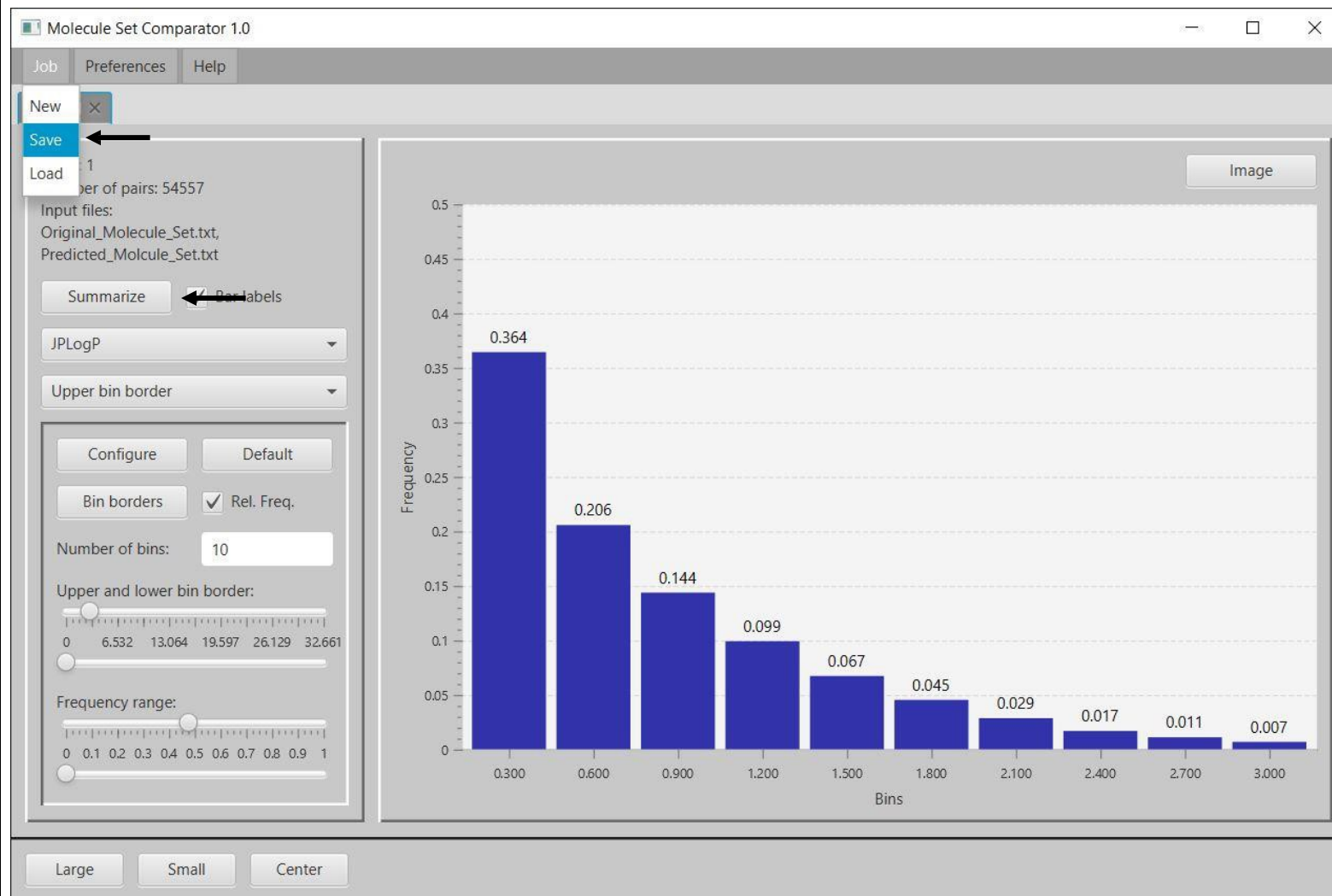


9

# Other descriptors

- The configured histogram of the molecular weight descriptor also exhibits a "nice" distribution without anomalies. The intervals are labelled with their upper border (so the first interval goes from 0 to 7.5)

- Roughly 36% of the predicted molecules have a similar molecular weight that only differs by about 0 to 7.5 daltons (see first bin)

- By clicking on the **Image** button, the histogram can be exported into various file formats for different purposes.

# Other descriptors

- The configured histogram of the JPLogP value can be seen on the right. The intervals are labelled with their upper border (so the first interval goes from 0.0 to 0.3)

- The histogram shows a distribution that is very similar to the one of the molecular weights. This connection could be worthy of a more thorough exploration with, e.g. the detail window

- Using the **Summarize** button, a summary document can be generated

- The processed job can be saved to be reloaded later or to be shared with others.

# Summary

- The goal here was to evaluate the capacity of a machine learning algorithm by comparing the predictions with the original molecule set.

- For the evaluation, four molecular descriptors were used (Tanimoto coefficient, atom count, JPLogP value, and molecular weight).

- Based on the histograms a quick assessment of the predictive power of the algorithm is possible. Here we could see that some molecules were predicted well, but the majority got predicted poorly. So there is room for improvement.

- Also, a new evaluation based on Tanimoto coefficients with other fingerprints is advised for a more thorough evaluation.

- Also, the tool can generate multiple outputs to share the results.