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

- Jan-Mathis Hein
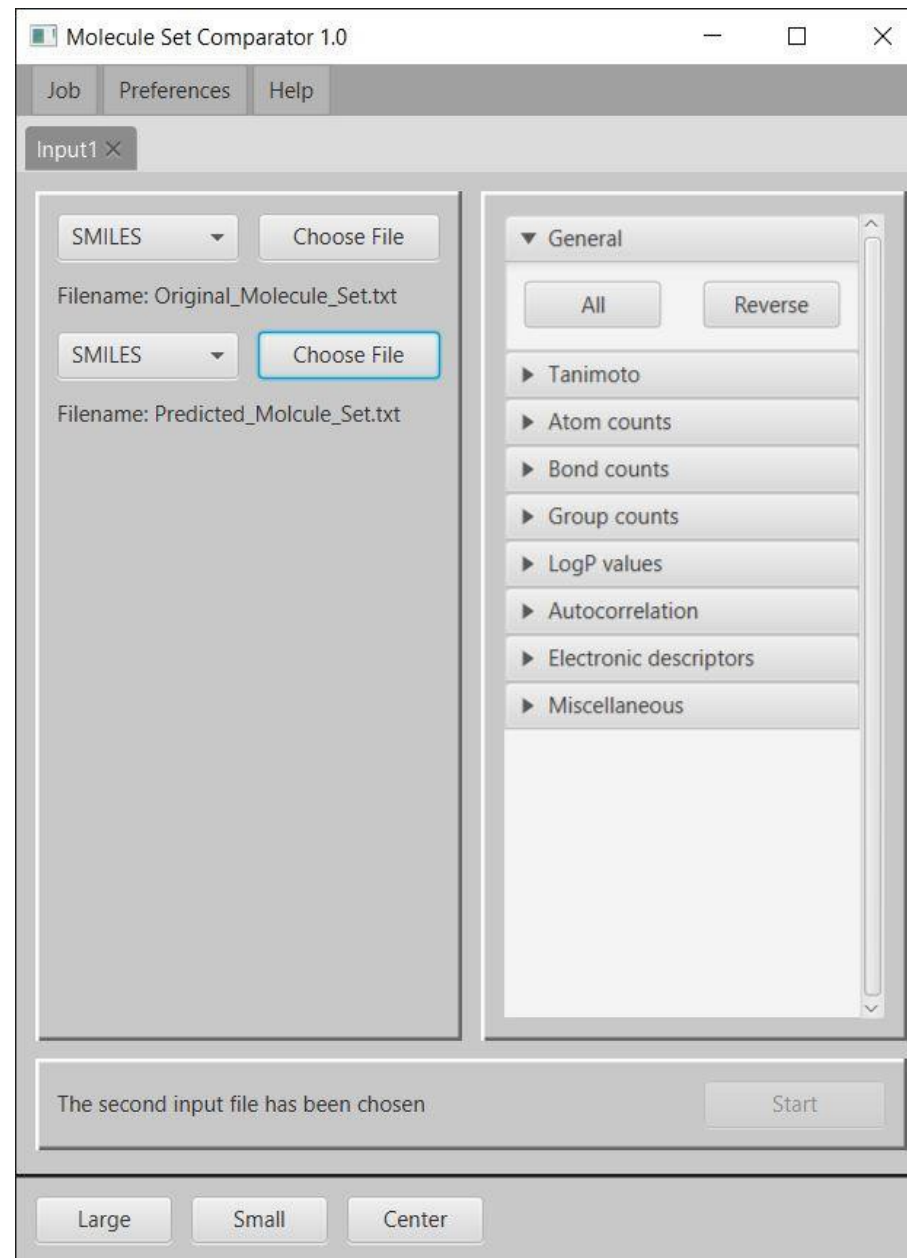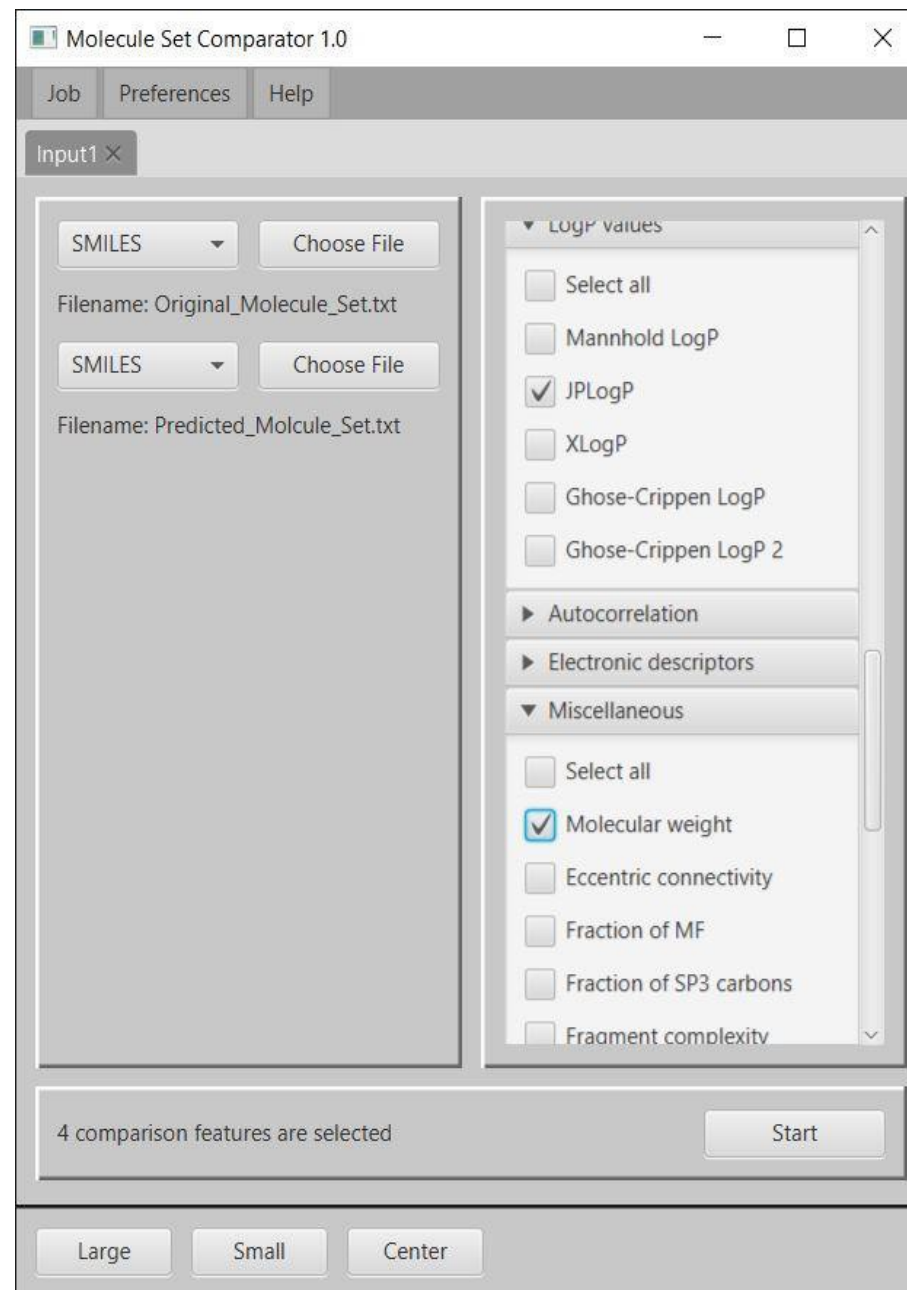
# Initializing the job

- In this tutorial we demonstrate how to use MSC to evaluate your machine learning predictions.

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

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

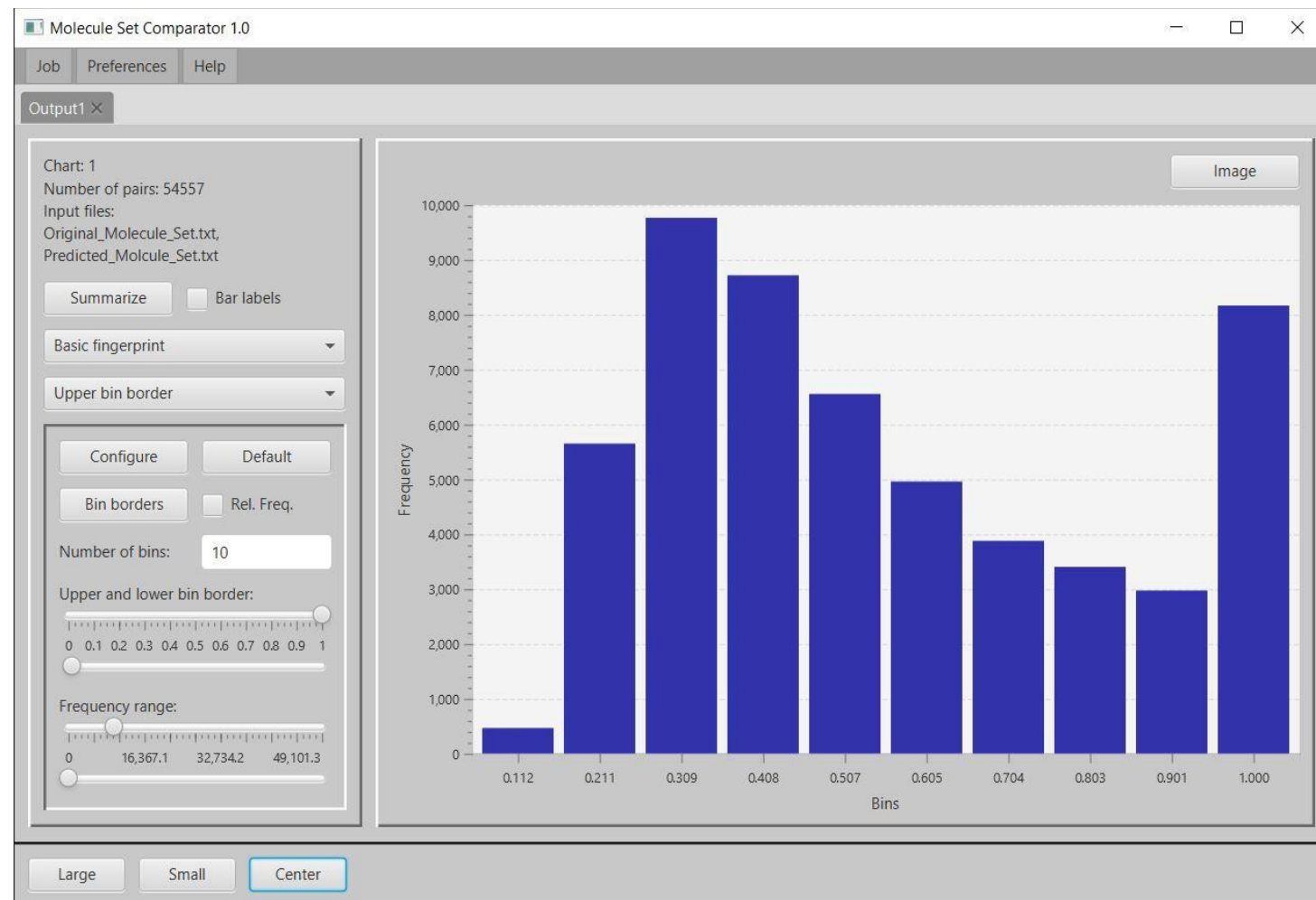- Here, both sets are encoded in the SMILES format.

# 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)
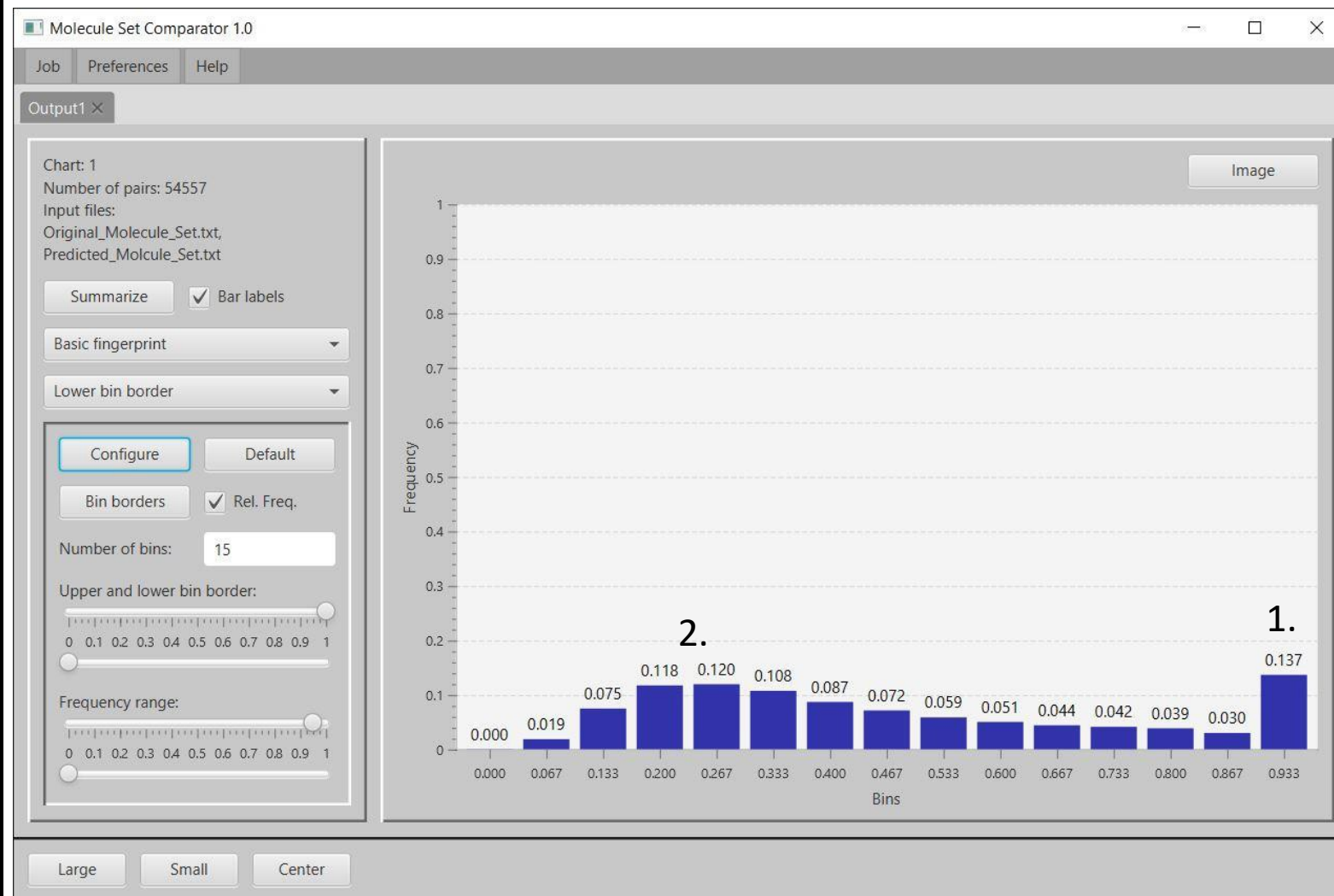
# Analysing the results

- 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 upper border of each interval is displayed.

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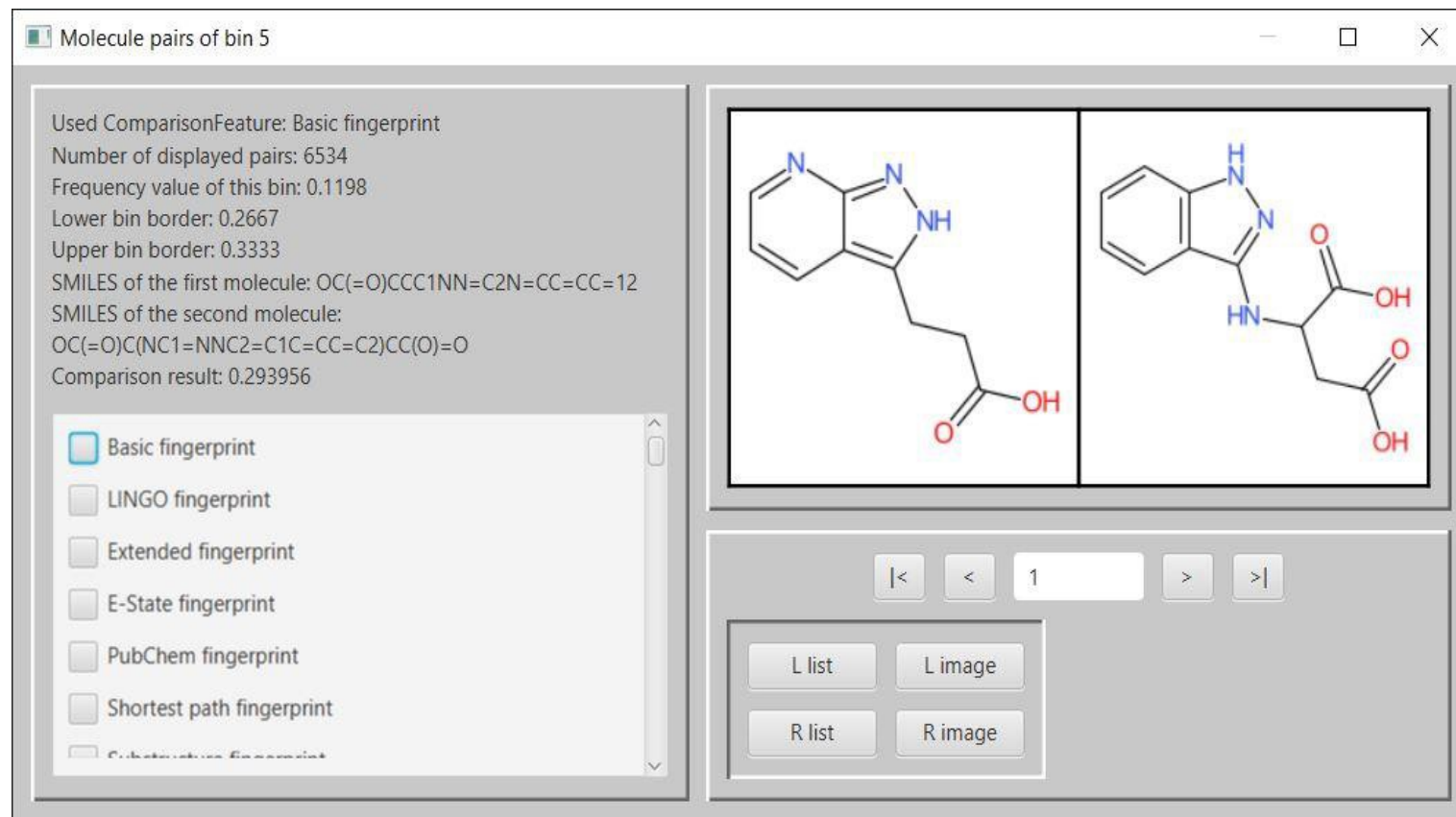
# Analysing the results

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

- Every bar is labelled with its frequency.

- The number of intervals and range of the y-axis range got increased.

- The x-axis labels now display the lower border of each interval.

- Now, the data can be easily analysed.

- Around 13.7% 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.

- On the right, the detail window for the fifth bin is displayed. Here one can browse through the list of input-output-pairs of this bin.

- It is immediately visible that the first molecule pair doesn't look that different. This intuitive rating contradicts the low Tanimoto coefficient of 0.29
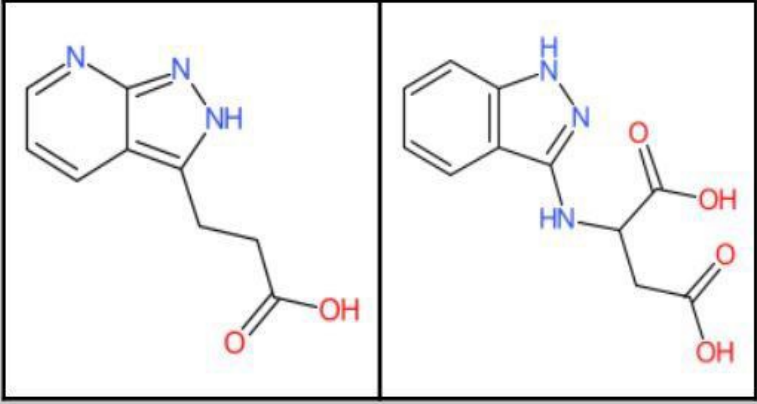
# The detail window

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

- For example, the PubChem fingerprint yields a Tanimoto value of 0.72 for the same molecule pairs which seem to be a lot more reasonable.

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

- To get a clear picture we can 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.

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

# Other descriptors

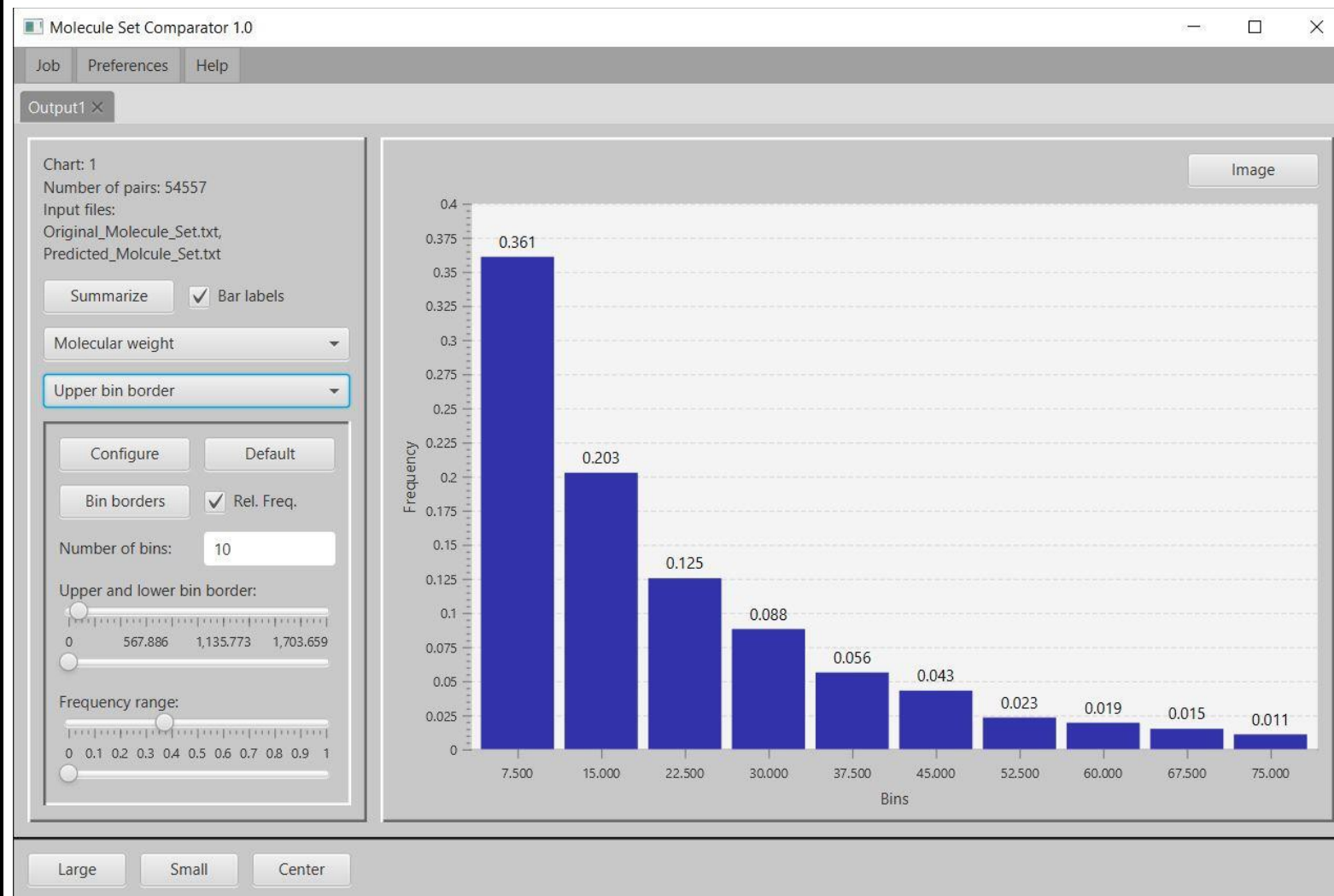- After configuring, the histogram conveys a lot more information.

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

- The histogram shows a better overall distribution except for the 4th bar.

- But with the help of the detail window, it can be concluded that, this anomaly is likely to be caused by the addition of one "C" atom which also leads to two additional "H" atoms.

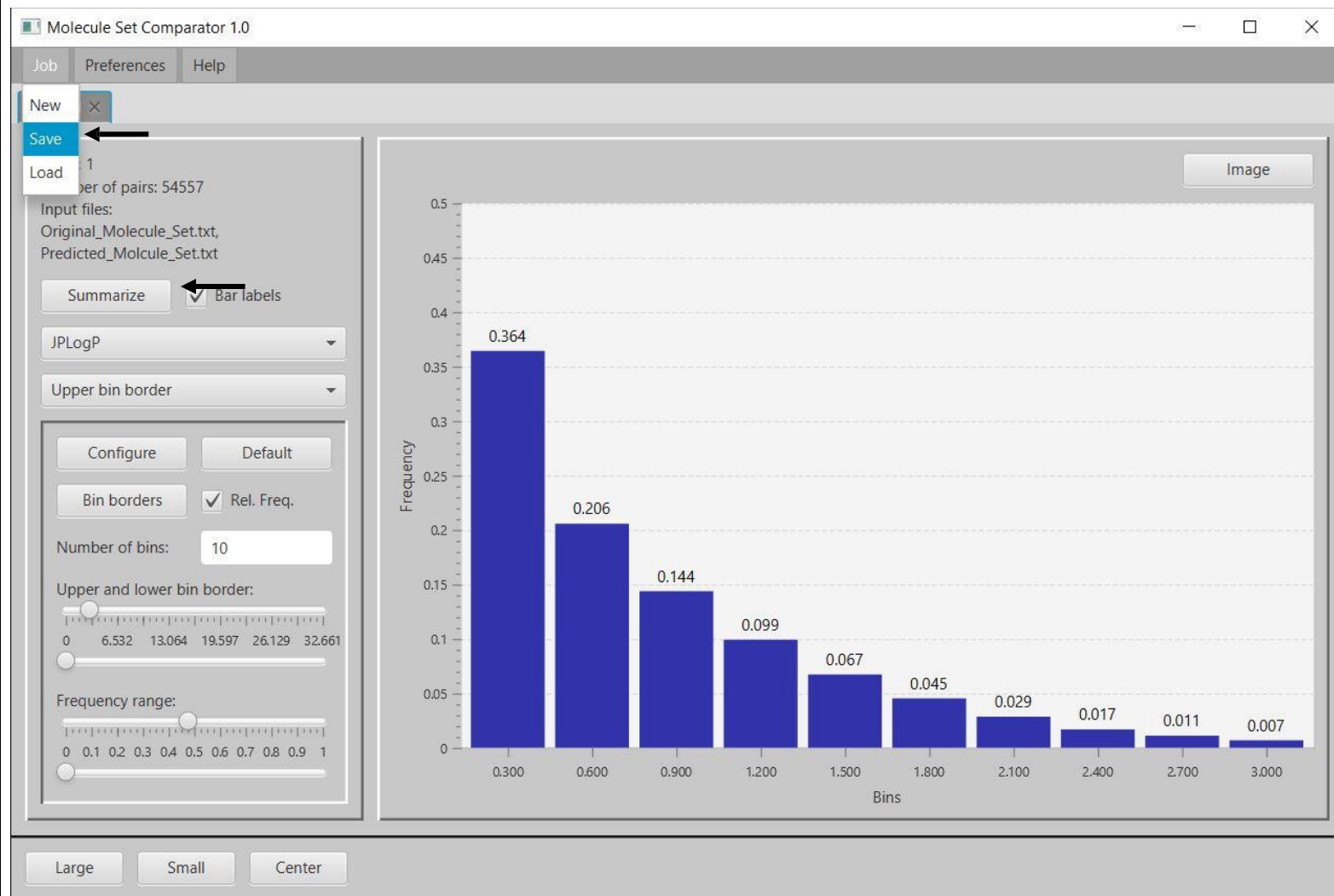# Other descriptors

- The configured histogram of the molecular weight descriptor also exhibits a better overall distribution without anomalies.

- Roughly, 36% of the predicted molecules have a similar molecular weight that only differs by about 0 to 7.5 daltons.

- Clicking on the **Image** button, the histogram can be exported into various file formats. This can later be used for any other purpose.

# Other descriptors

- The configured histogram of the JPLogP value can be seen on the right

- It shows a distribution that is very much similar to the one with 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 whole processed job can be saved to check these results later or to share them 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. The results generated using the JPLogP and molecular weight descriptors can give a better insight into the predictions as well.

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