

Worked Example:

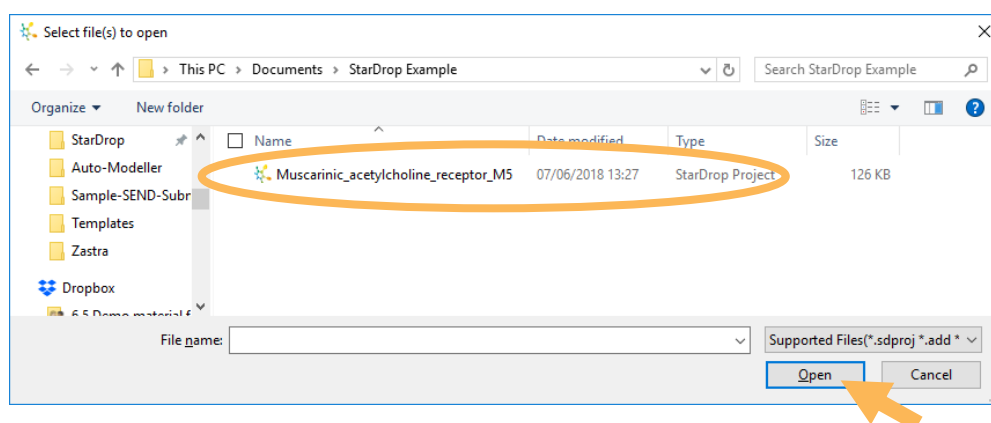
Automatic QSAR Model Building and Validation

In this example, we will explore the application of StarDrop's Auto-Modeller to build a QSAR model of potency against the Muscarinic Acetylcholine M5 receptor, based on a set of public domain K_i data obtained from the ChEMBL database (<https://www.ebi.ac.uk/chembl/>). The resulting model will be applied to an additional set of compounds to predict their properties and visualise the structure activity relationship.

Step-by-step instructions for all the features you will need to use in StarDrop are provided, along with screenshots and examples of the output you are likely to generate. If you have any questions, please feel free to contact stardrop-support@optibrium.com.

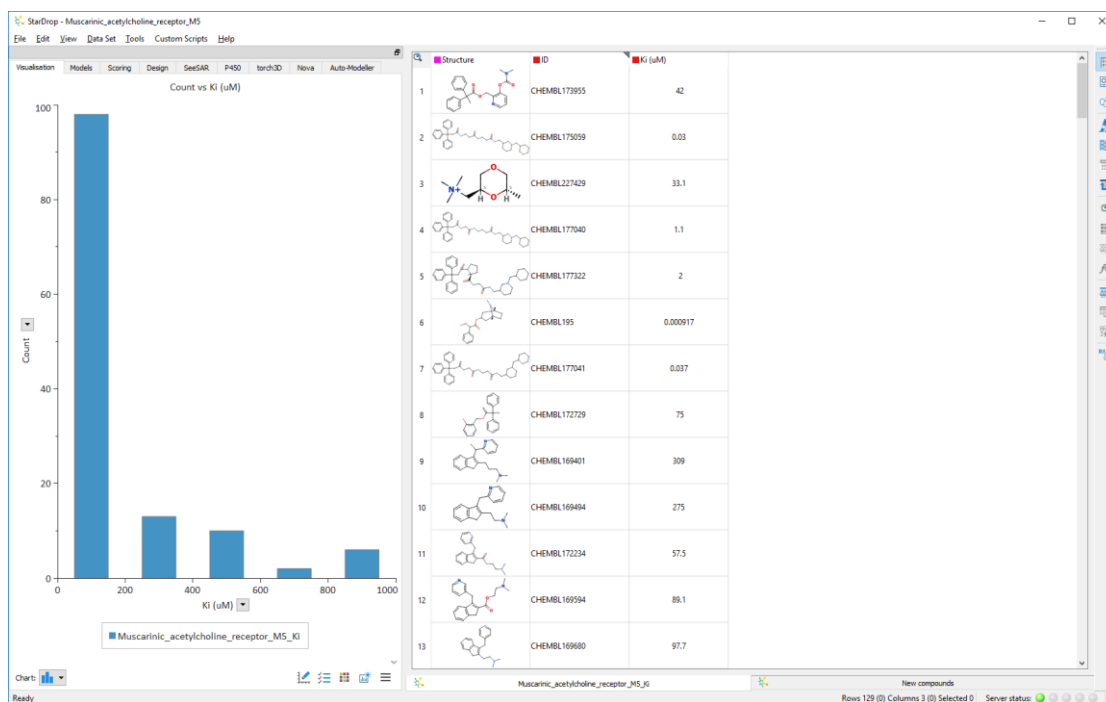
Exercise

- Open the StarDrop Project file **Muscarinic_acetylcholine_receptor_M5.sdproj** by selecting **Open** from **File menu**.



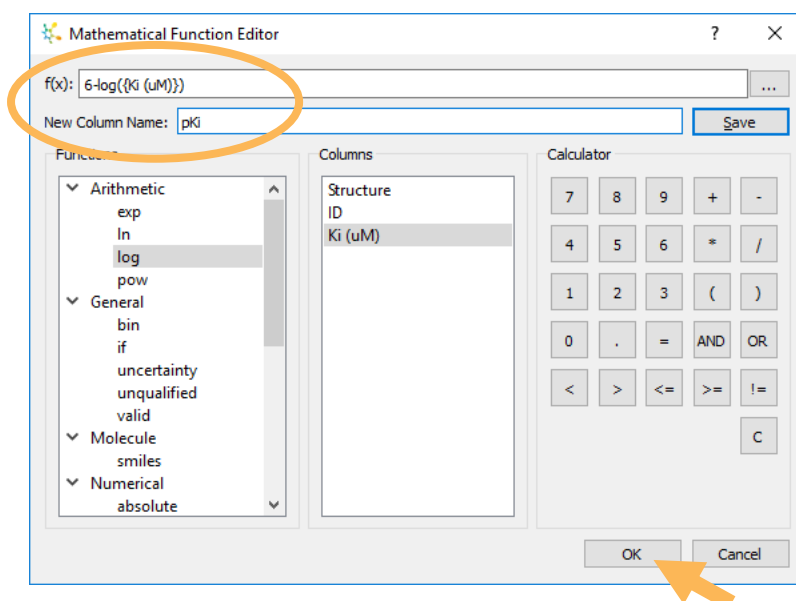
Optibrium™, StarDrop™, Card View™, Nova™ Glowing Molecule™ and Auto-Modeller™ are trademarks of Optibrium Ltd.

© 2020 Optibrium Ltd.

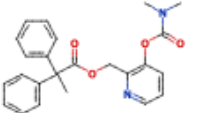

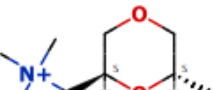



The first data set displayed contains 129 compounds, with measured activity against the muscarinic acetylcholine receptor M5. These data are K_i values in μM . However, to build a good model the data should be converted to logged units. Logged units provide a more even distribution of values to model and a better correlation with the compound descriptors used to build the model. Therefore, we will use the **Mathematical Function Editor** in StarDrop to generate $\text{p}K_i$ values.

- Select the $f(x)$ button on the toolbar to open the **Mathematical Function Editor**
- In the $f(x)$ field, enter the equation “ $6 - \log(\{K_i (\mu\text{M})\})$ ”. This can be easily achieved by pointing and clicking in the editor (or by copying and pasting without the quotes). Enter the name of the new column, **$\text{p}K_i$** , in the **New Column Name** field and click **OK**.



The new column containing pK_i values will appear in the data set. Now we're ready to build a model of this data.

	Structure	ID	Ki (uM)	pKi
1		CHEMBL173955	42	4.38
2		CHEMBL175059	0.03	7.52
3		CHEMBL227429	33.1	4.48

- Change to the **Auto-Modeler** area on the left and click on the  button to begin a new modelling session

	Structure	ID	Ki (uM)	pKi
1	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL173955	42	4.38
2	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL173059	0.03	7.52
3	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL227429	33.1	4.48
4	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL177040	1.1	5.96
5	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL177332	2	5.7
6	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL195	0.000917	9.04
7	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL177041	0.037	7.43
8	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL172729	75	4.12
9	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL169401	309	3.51
10	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL169494	275	3.56
11	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL172234	57.5	4.24
12	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL169594	86.1	4.05
13	<chem>CC1=CC=C(C=C1)C2=CC=CC=C2</chem>	CHEMBL169680	97.7	4.01

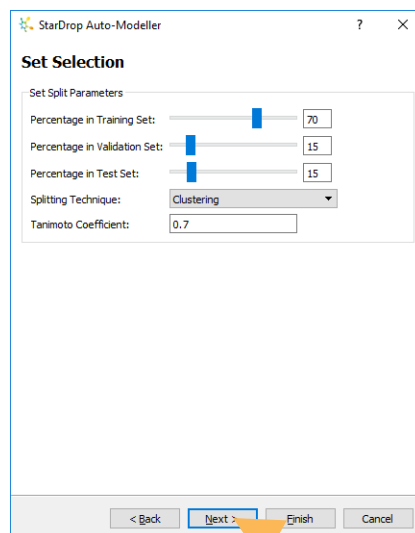
This will open the Auto-Modeller wizard.

The first page enables you to choose the type of model to build: continuous (i.e. numerical) or category (i.e. classification). You can also choose whether to allow StarDrop's Auto-Modeller to split the data into Training, Validation and Test sets automatically or to provide these separate sets yourself. Finally, you can confirm the data set to model and the columns containing the property values to model, the compound structures and compound identifiers. In this case, we will use the default settings, but we need to select the correct column to model.

- Choose the **pKi** column from the **Value Column** drop-down menu and click **Next**

The next page enables you to configure the parameters for the automatic selection of Training, Validation and Test sets. In this case we will use the default set split parameters.

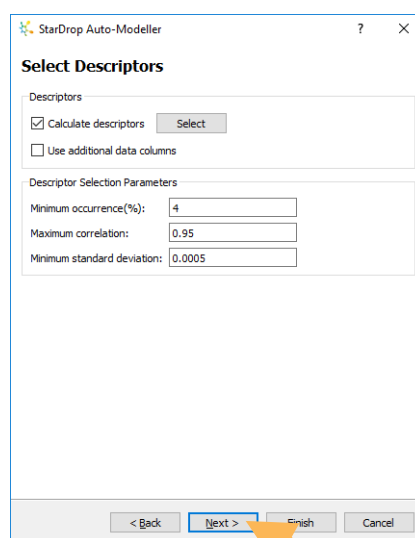
- Click **Next**



The 'Set Selection' dialog in StarDrop Auto-Modeller shows default split parameters. It includes sliders and input fields for 'Percentage in Training Set' (70), 'Percentage in Validation Set' (15), and 'Percentage in Test Set' (15). The 'Splitting Technique' is set to 'Clustering' and the 'Tanimoto Coefficient' is 0.7. An orange arrow points to the 'Next >' button at the bottom.

The third page enables you to select the descriptors to use. The built-in library of whole molecule and 2D descriptors will be used by default. More details are available by clicking the **Select** button and new descriptors can be imported as SMARTS. Also, additional columns in the data set can be used as descriptors. Finally, the parameters for selection of descriptors can be defined. Again, we will use the default settings.

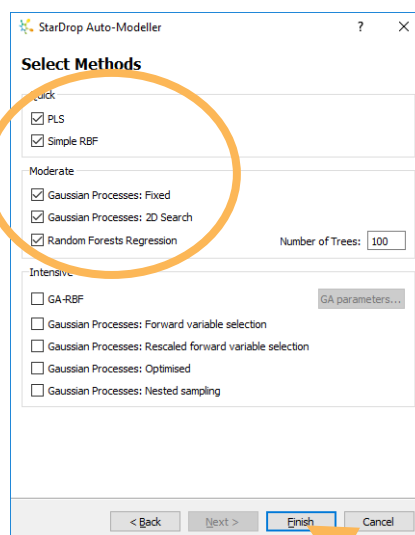
- Click **Next**



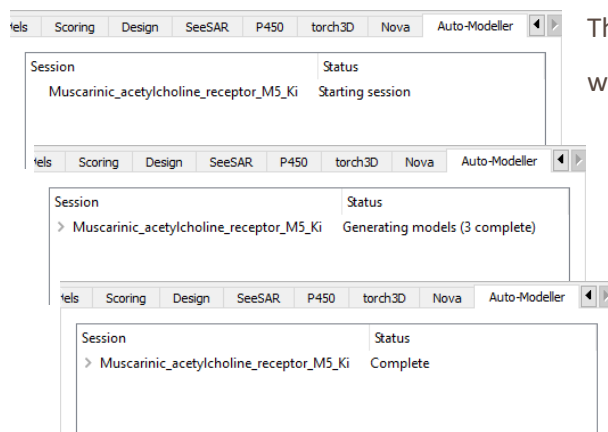
The 'Select Descriptors' dialog in StarDrop Auto-Modeller shows default descriptor selection parameters. It includes checkboxes for 'Calculate descriptors' (checked) and 'Use additional data columns' (unchecked). The 'Descriptor Selection Parameters' section has input fields for 'Minimum occurrence(%)' (4), 'Maximum correlation' (0.95), and 'Minimum standard deviation' (0.0005). An orange arrow points to the 'Next >' button at the bottom.

The final page of the Auto-Modeller wizard enables you to select the modelling methods to apply. The methods are categorised by their computational cost and the methods selected by default will depend on the size of the data set.

- We would recommend unselecting the **Intensive** methods for this quick example, as shown to the right.
- Click **Finish** to begin the modelling session

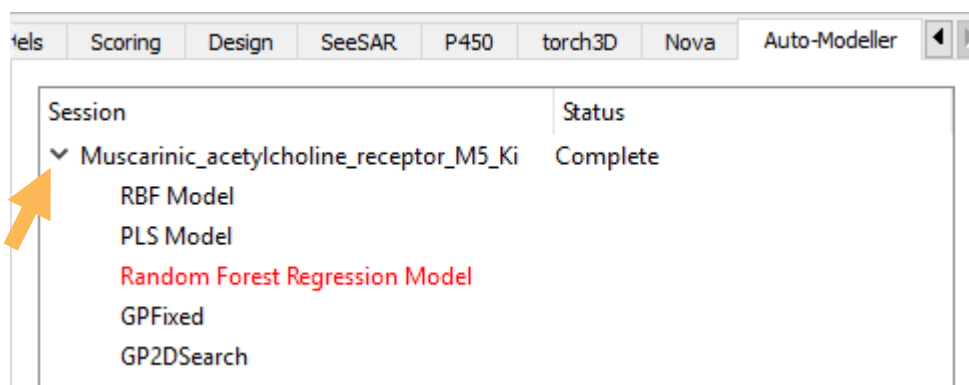


The 'Select Methods' dialog in StarDrop Auto-Modeller shows modelling methods categorized by computational cost. The 'Quick' category has 'PLS' and 'Simple RBF' selected. The 'Moderate' category has 'Gaussian Processes: Fixed', 'Gaussian Processes: 2D Search', and 'Random Forests Regression' selected, with 'Number of Trees' set to 100. The 'Intensive' category has 'GA-RBF' and several 'Gaussian Processes' options. An orange circle highlights the 'Quick' and 'Moderate' categories, and an orange arrow points to the 'Finish' button at the bottom.



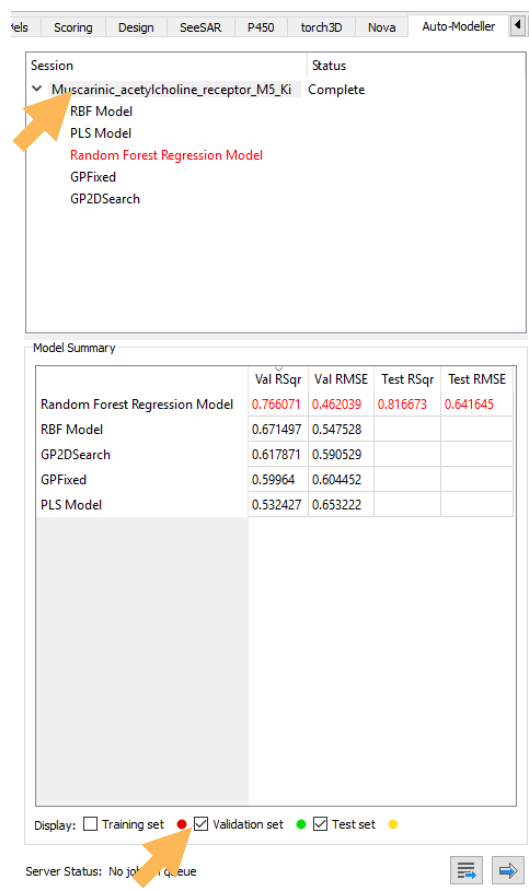
The top section of the **Auto-Modeller** area will provide a running update on the progress of your modelling session. If there are no other modelling sessions ahead of yours in the queue on the server, this should only take a few minutes (you can check the status of the server at the bottom of the **Auto-Modeller** area).

- When the modelling session is complete, click the arrow next to the session in the **Auto-Modeller** area to open up the list of models generated.



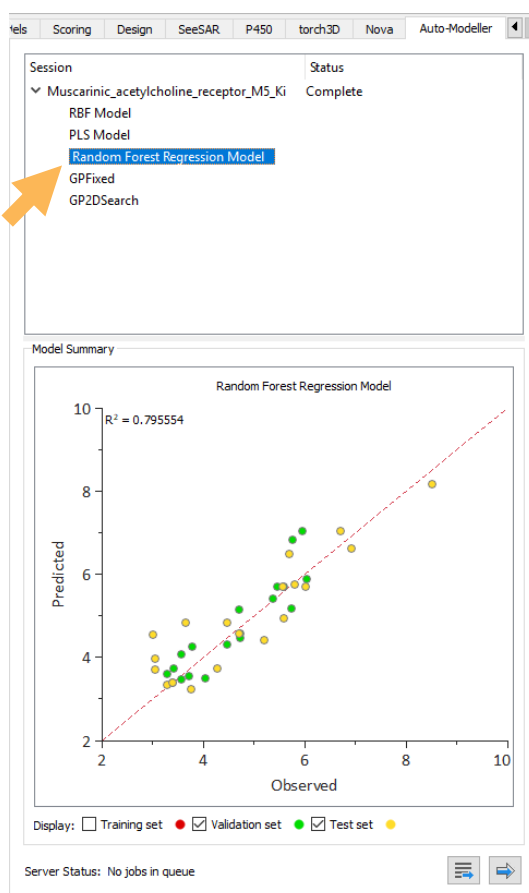
The model with the best result on the Validation set will be highlighted in red. Please note that the specific results you see may be slightly different to the examples shown here due to a random element in the assignment of compounds to the Training, Validation and Test sets.

- Select the modelling session to see a table summarising the results of the different models. This will show the result of the best model on the Test set.
- Tick the **Validation set** box at the bottom of the **Auto-Modeller** area to see the results for all of the models on the Validation set.

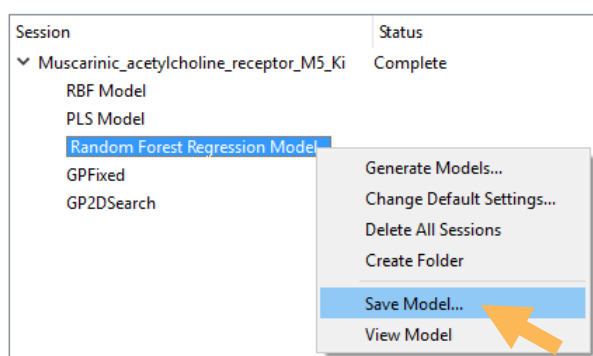


You can see that the Auto-Modeller has built a model with each of the modelling methods and compared the performance of these for the Validation set to identify the most predictive model. This best model is then further validated using the external Test set. A robust model should have good performance on both the Validation and Test sets.

- Select a model to see a plot of the results for the Validation and Test sets and confirm that the model is producing reliable predictions.
- Hover the mouse pointer over a data point to see the corresponding structure.

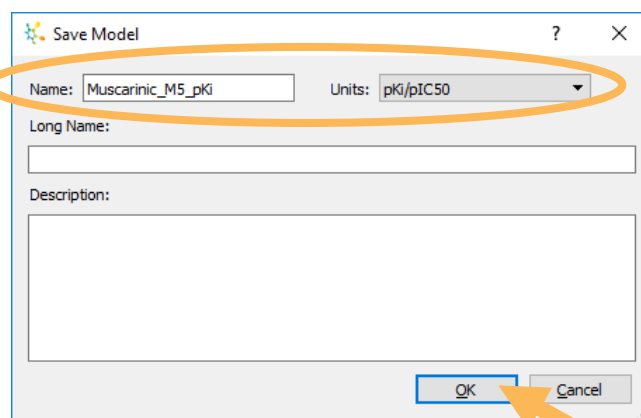


If you are happy with the results for a model, you can save it so that you can apply it to new compounds to predict the potency and visualise the structure-activity relationship. In this case we will save the random forest regression model.

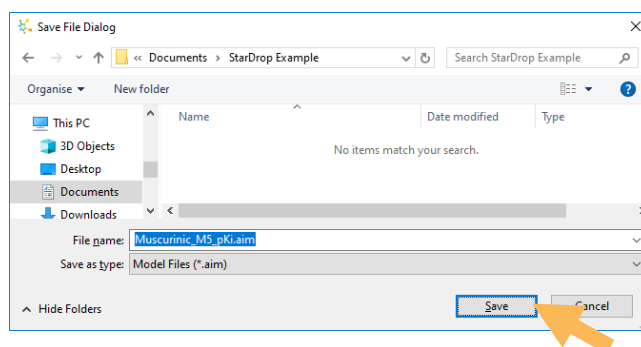


- Right click on the **Random Forest Regression Model** under the modelling session in the **Auto-Modeller** area and select the **Save Model** option.

- In the **Save Model** dialogue enter an appropriate name, set the units to **pKi/pIC50** and click **OK** (if you wish you can also enter more information in the **Long Name** and **Description** boxes).



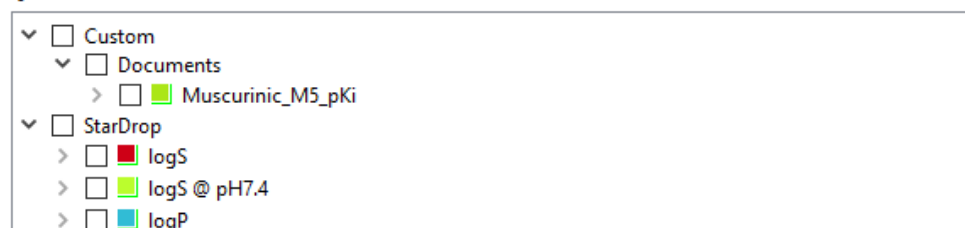
- Finally, navigate to a convenient directory and click **Save** to save the model.



Note: The resulting file can be shared with any other StarDrop user who can load and use the model in their copy of StarDrop.

- Switch to the **Models** area and you will see that the model has appeared under the corresponding directory name, ready to run on new compounds.

QSAR Models:




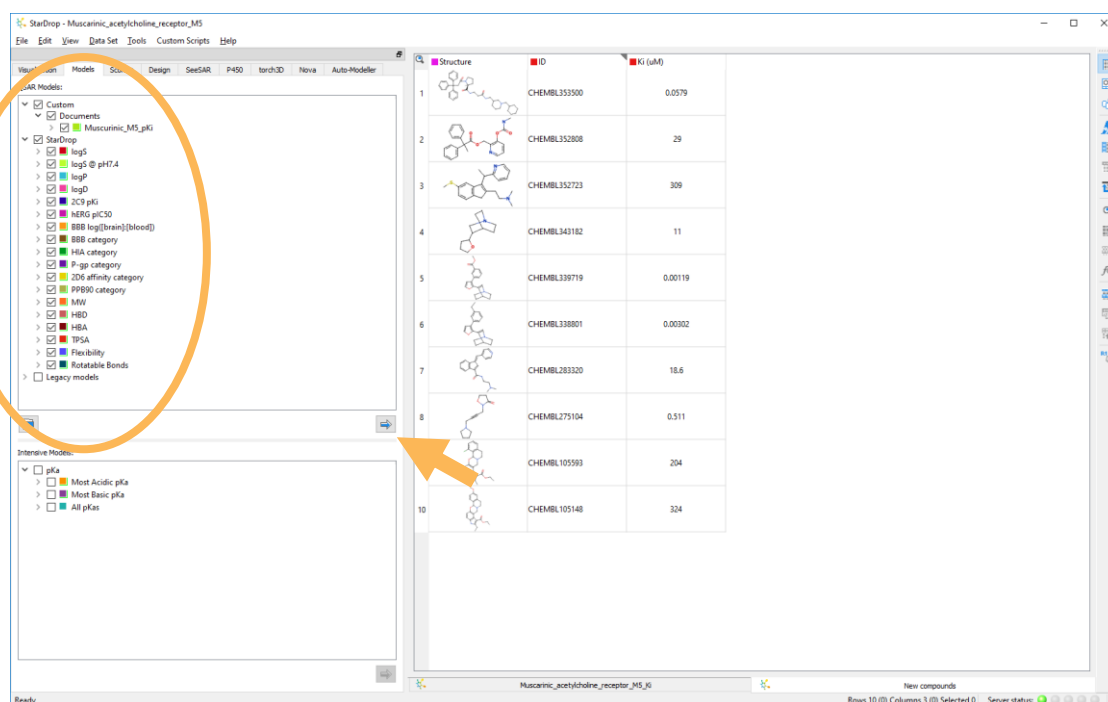
We'll illustrate the application of this model to an additional 10 compounds that were not included in the data set used to build and validate the model.

- Change to the **New compounds** data set by clicking on the tab at the bottom of the data set.

The screenshot shows the StarDrop software interface for the Muscarinic_acetylcholine_receptor_M5 model. The 'Visualisation' panel on the left lists various models under the 'Muscarinic_M5_pKi' category. The 'Structure' panel on the right displays a table of 10 compounds with their IDs and Ki values. An orange arrow points to the 'New compounds' tab at the bottom right of the interface.

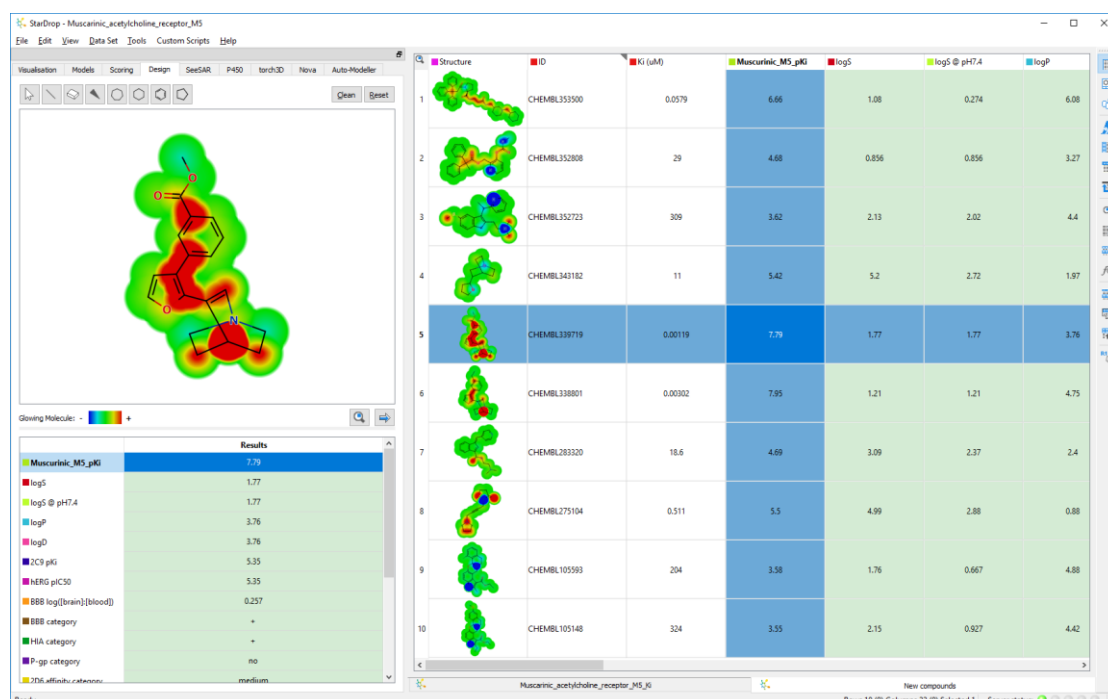
Structure	ID	Ki (uM)
1	CHEMBL353500	0.0579
2	CHEMBL352808	29
3	CHEMBL352723	309
4	CHEMBL343182	11
5	CHEMBL359719	0.00119
6	CHEMBL338801	0.00902
7	CHEMBL283320	18.6
8	CHEMBL275104	0.511
9	CHEMBL105593	204
10	CHEMBL105148	304

- In the **Models** area, select the new model we have built along with any other models you would like to run, by ticking the boxes next to the models, and click the  button.



The new model can be used in the same way as any other model in StarDrop. For example, selecting the column header will display the Glowing Molecule visualisation for each compound, showing the structure-activity relationship captured by the model we have built.

Changing to the **Design** area and selecting a row in the data set will enable you to explore optimisation strategies, guided by the Glowing Molecule.



This has been a quick example of the application of StarDrop's Auto-Modeller. There are, of course, additional features enabling expert modellers to control the parameters of the

model building process and explore the detailed results for each model. For more information or to arrange a comprehensive demo, please contact stardrop-support@optibrium.com.