



## Worked Example:

### Generating New Optimisation Ideas Using Matched Series Analysis

The objective in this worked example is to identify new derivatives that are likely to improve activity at their target, given the SAR already generated on a project. This example uses a publicly available set of Human Histamine H1  $K_i$  data and searches the ChEMBL  $pIC_{50}$  knowledge base (generated by NextMove Software) to find matched series that identify new substitutions with a high likelihood of having improved binding.

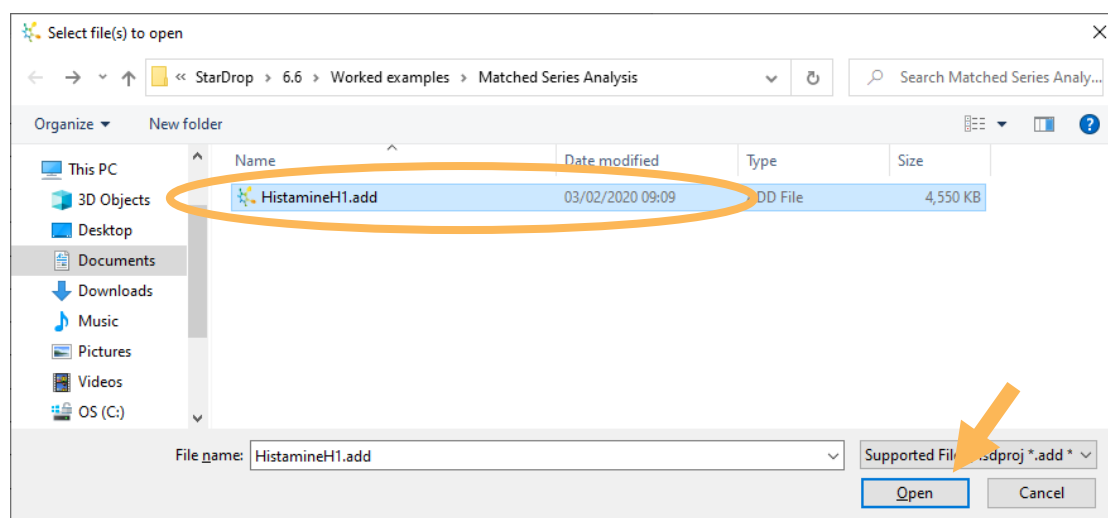
A matched series is a series of compounds that are identical except for different substituents at a single point (see Section 10.3 of the StarDrop Reference Guide for more details). The suggestions are derived from comparing matched series in the input data with those in a knowledge base, which are measured across diverse target proteins. The suggestions are based on the premise that a matched series with similar activity order in the input data set and the knowledge base implies that those groups occupy a similar binding environment created by their target proteins. Given a similar binding environment, groups that have been shown to be better binders within the knowledge base, have a strong likelihood of being better binders to the target behind the input data set, in this case Histamine H1.

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](mailto:stardrop-support@optibrium.com).



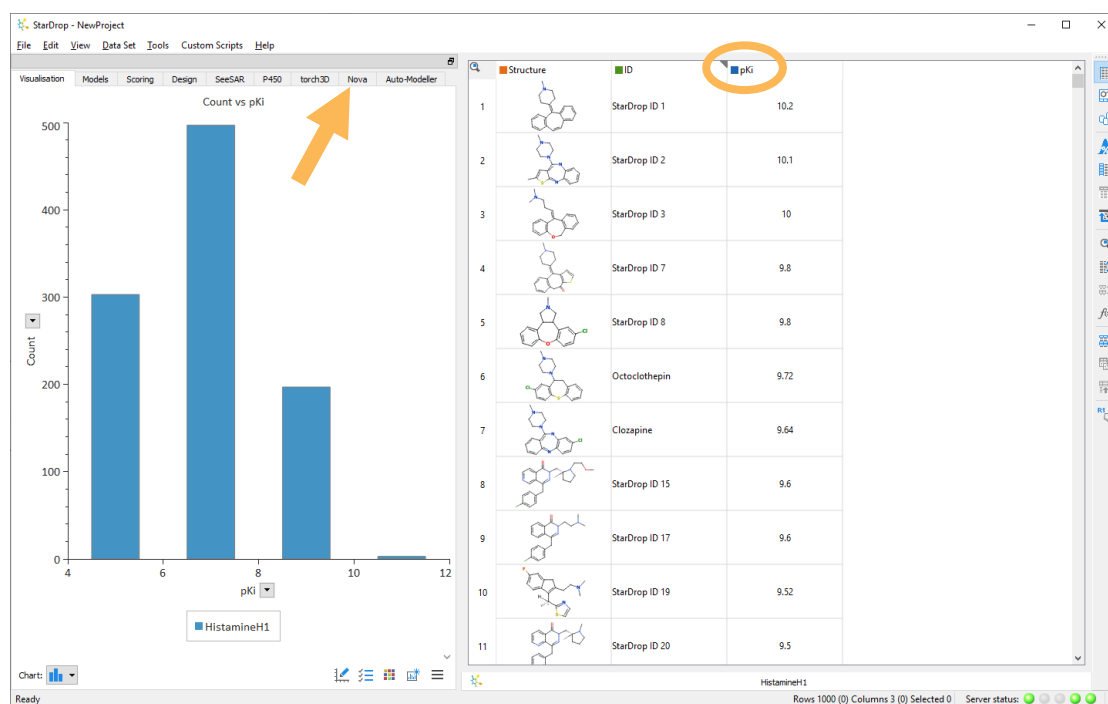
Optibrium™, StarDrop™, Card View™, Nova™, Glowing Molecule™ and Auto-Modeller™ are trademarks of Optibrium Ltd.  
Matsy™ is a trademark of NextMove Software Ltd.


## Exercise



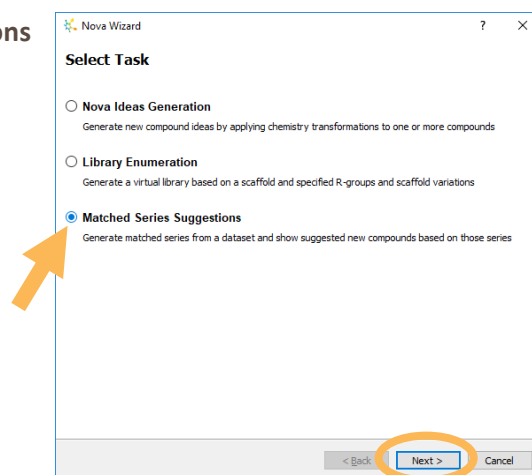
Open the file **HistamineH1.add** by selecting **Open** from the **File** menu.

- You will see a spreadsheet containing 1000 structures and their measured affinities for Human Histamine H1 (in the column labelled **pKi**).

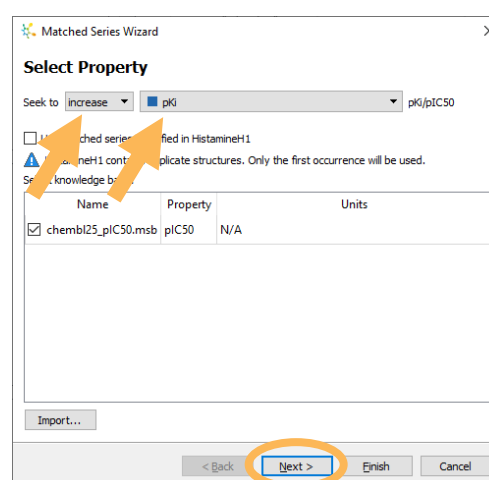


- To start the matched series analysis, click on the **Nova** tab and then at the bottom of the **Nova** area click the  button.

- Select the **Matched Series Suggestions** option and click **Next**



- In the dialogue box that appears, you can specify the column containing the property you wish to improve. In this case, the column we are interested in, **pKi**, is already chosen and we want to find suggestions that **increase** this value, so this default option is also correct.
- Select **Next** to continue.



At this point you can change the limitations placed on the suggestions returned. In Matsy™ the support for a suggestion comes from the number of times it has been seen; the more frequent the occurrence of the order of the input series, the more likely it is that the suggestion will be an improvement. Hence, to find many examples in the ChEMBL knowledge base, the compared series are generally short.

The default options are to match a series of 3 derivatives and that series should have been seen at least 20 times in the knowledge base and these are acceptable for this data set.

With SAR transfer, the support for a suggestion comes from a long series of derivatives that shows a consistent trend with that seen in the input data set. This example data set is too small to have matched series with the default minimum number of derivatives (8), so for this example we will decrease this limit.

- Click on the **Minimum length of matched series** box in the **SAR transfer** section and change the value to **7**.
- Click the **Next** button to continue.

Here you can give the output data set of suggestions a different name and control what else is reported in the output.

- Check all the boxes for the **Matsy** output as shown above and select **Next**
- Here you can choose structural filters to exclude certain chemical groups. In this case we will use the defaults, so click **Finish** to begin the matched series analysis.

**Matched Series Wizard**

**Select Method**

☒ Calculate Matsy suggestions

Minimum number of matched series: 20

Minimum length of matched series: 3

☒ Calculate SAR transfer suggestions

Minimum series correlation: 0.7

Minimum length of matched series: 7

< Back **Next >** Finish Cancel

**Matched Series Wizard**

**Output Data Set**

Name: Chymotrypsin\_pK\_i\_suggestions

**Matsy**

☒ Percent that improve

☒ Total number of observations

☒ Enrichment

**SAR Transfer**

☒ Maximum correlation

☐ Number of series with improving SAR transfer

☐ Show in Card View

< Back **Next >** Finish Cancel

The suggestions are returned in a table with the Matsy based suggestions first, followed by the SAR transfer suggestions. The Matsy suggestions are ordered by the **% that improve** column and the SAR transfer are ordered by the **maximum correlation** column. When a row is selected in the data set the suggestion is displayed in the Nova area and the supporting evidence is shown in a table below.

**Suggestion:**

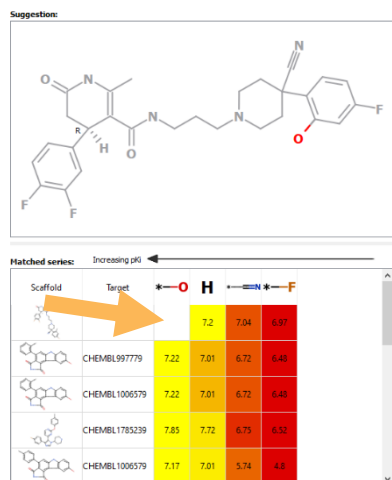
**Matched series:**

Scaffold	Target	*-O	H	*-N	F
		7.22	7.04	6.97	
CHEMBL997779		7.22	7.01	6.72	6.48
CHEMBL1006579		7.22	7.01	6.72	6.48
CHEMBL1785239		7.85	7.72	6.75	6.52
CHEMBL1006579		7.17	7.01	5.74	4.8

**Table:**

	Structure	R-Group	Scaffold	% that improve	Enrichment
1				51.9	1.34
2				44.1	1.44
3				42.3	1.56
4				40.7	2.64
5				39.7	1.32
6				39.1	1.49
7				34.8	1.05
8				33.3	1.56
9				32.7	1.21
10				32.1	0.783
11				31.3	0.63

The SAR data from the input data set is in the first row (which is why the target and first substituent columns are empty) and the SAR data are ordered with the least active/desirable

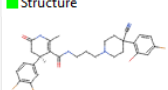
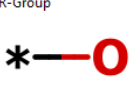
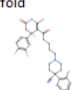


on the right to the most active/desirable on the left (as indicated by the colour coding in the table cells).

The first row of the data set shows one of the suggestions that is most likely to improve the pK<sub>i</sub> which is the creation of the hydroxy derivative. This suggestion is based on the order of activity seen for the fluoro, cyano, and unsubstituted derivatives, at that position on the displayed scaffold, seen in the input data set.

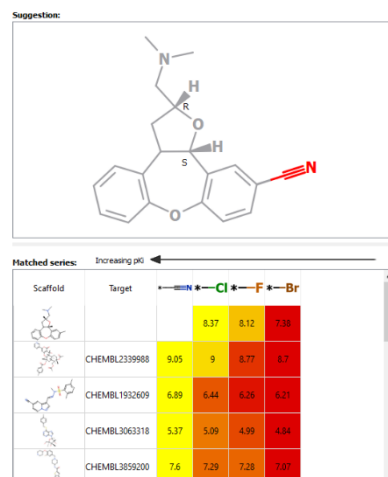
The supporting evidence for the suggestion comes from a variety of target sources and scaffolds and more information for each entry can be obtained by either clicking on the target name in the target column (which will bring up the ChEMBL web page for that target) or by hovering the mouse over the scaffold image in the table to give an enlarged view.

The first of the matching series from the ChEMBL knowledgebase is for target CHEMBL997779 (which is human WEE1 kinase) where the series also occurs at the ortho position, matching the scaffold in our series. Note that further examples are also seen in the same order for the meta and para positions, which may relate to influencing the aryl ring electronics rather than any direct effects. This pattern has been seen many times and nearly 52% of the 27 observations have shown an increase in activity.

Structure	R-Group	Scaffold	% that improve	Enrichment	Observations
			51.9	1.34	27

- Select row 13 in the data set.

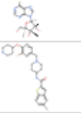
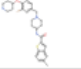
Here we have an interesting example because it does not follow the traditional order of halogens by molecular weight. In addition, the cyano group is not a group one might typically consider next in the series, yet there are 49 observations of this series in the knowledgebase where this results in an increase in activity nearly 31% of the time.



**Note:** You can view the

examples where the order

does not match the input series for any of the suggestions by ticking the **Show negative** option below the table of examples.

	CHEMBL3063318	5.37	5.09	4.99	4.84
	CHEMBL3859200	7.6	7.29	7.28	7.07

☒ Show heatmap

☒ Colour by row

☒ Show negative

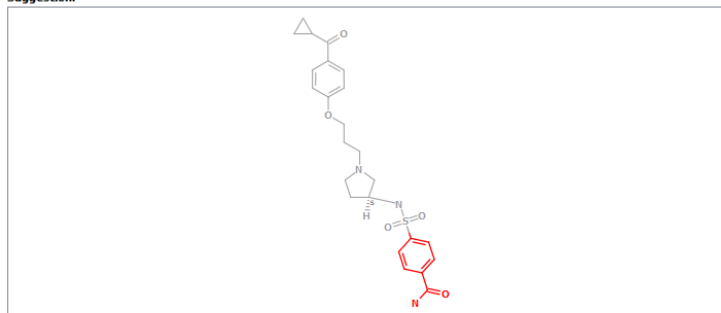
Colours...

- To see the SAR transfer suggestions, **right-click** on the **Max. Correlation** column in the data set and choose **Descending** from the **Sort** menu.

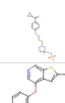
In the top suggestion you can see that the series of aryl derivatives from the input data correlates very well with the activities of derivatives at MAP3K8. **Note:** You may need to make your Nova area wider to see the complete series by dragging with the mouse.

Max. Correlation	Series	pKi
Delete		
Insert...		
Duplicate		
Sort		
Sort by Confidence		
Edit...		
Copy		
Add Data Set to Fragment Library...		

Suggestion:



Matched series:

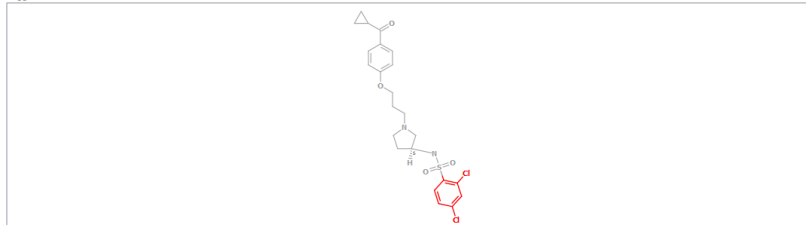
Scaffold	Target	Correlation	5.6	5.44	5.36	5.35	5.04	5.04	4.7
	CHEMBL1005354	0.991	8	7.05	6.92	6.89	6.82	6.77	6.51

For this suggestion there is only one example, whereas some suggestions are based on multiple long series.

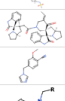
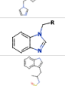
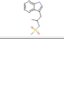

- Select row 4 in the data set.

This series is very long and has multiple matches that correlate well. This is a good example where it would probably be difficult to determine

Suggestion:



Matched series:

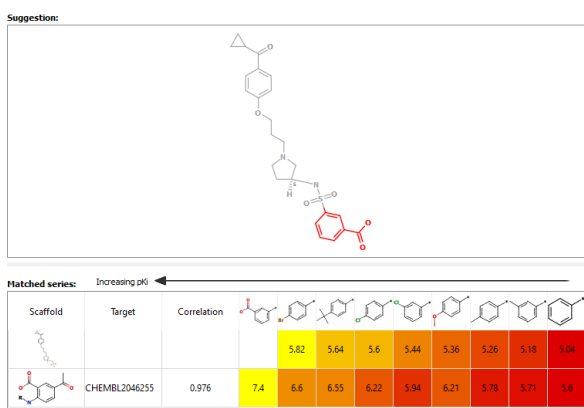
Scaffold	Target	Correlation	5.82	5.64	5.6	5.44	5.36	5.35	5.26	5.23	5.18	5.05	5.04	5.04	5.03	4.8
	CHEMBL2345676	0.991	5.16		5.1			4.87	4.71				4.71	4.63	4.61	4.56
	CHEMBL685883	0.786	6.55		6.04	5.49	5.46	5.96		5.36		4.8	5.43			
	CHEMBL1947960	0.919	4.66	4.56	4.57	4.07	4.42		4.05				4	4		
	CHEMBL890119	0.839	6.32		6.13	5.68	5.39	5	5.39	5.24		5	5	5		

such evidence for the suggestion based upon database searching, particularly given the various differences, and gaps, in the example series.

- Select row 6 in the data set.

Here we can see an example where the suggested compound is a benzoic acid, which will make the suggested scaffold zwitterionic. We may consider this suggestion surprising based on the fact that the reference sequence is all hydrophobic. While interesting and not

the obvious next derivative to make, this may come from the exemplar structure already being a benzoic acid and the target CHEMBL2046255 being aldo-keto-reductase (i.e. this may be a special case specific to that protein).



This worked example has shown how matched series analysis can generate suggestions for novel derivatives to improve the binding at your target, based on the data already generated within your own project. For any suggestions that appear interesting, it is important to consider the evidence provided by the example series to determine how this may translate to your own project. The applicability and suitability of suggestions also relies on many other compound properties so the data set of suggestions can be further prioritised in StarDrop, for example using predictive models of physicochemical and ADME properties or target activity and Probabilistic Scoring for multi-parameter optimisation.

Examples of the use of these methods can be found in further worked examples, but if you would like to see a demonstration please contact [stardrop-support@optibrium.com](mailto:stardrop-support@optibrium.com).