# Cheminformatics Machine Learning for Homogeneous Catalysis

**Schrödinger** | **Materials**

# Cheminformatics Machine Learning for Homogeneous Catalysis

**Topics:** Catalysis and Reactive Systems, Semiconductors, Energy Capture and Storage, Metals Alloys and Ceramics, Informatics and Team Collaboration

**Methodology:** Machine Learning

**Products Used:** MS Maestro, AutoQSAR, MS Informatics

This tutorial is written for use with a 3-button mouse with a scroll wheel.

Words found in the **Glossary of Terms** are shown like this: Workspace

**Abstract:**

In this tutorial, we will learn to develop and use a machine learning model to predict reaction rate constants for iridium catalysts.

**Tutorial Content**
1. Introduction
2. Creating Projects and Importing Structures
3. Building a Machine Learning Model Using DeepAutoQSAR
4. Viewing the Machine Learning Model and Predicting
5. Conclusion and References
6. Glossary of Terms

# 1. Introduction

Discovering new catalysts for improved reactivity or selectivity is challenging because of the large number of laborious experiments or stepwise quantum mechanical calculations necessary to explore the catalyst design space. Alternative to these approaches, employing machine learning (ML) for catalyst discovery and design is a promising avenue to rapidly screen catalysts for enhanced properties (see **References** for recent literature examples).

A useful ML tool is Quantitative Structure-Activity Relationships (QSAR), which can efficiently predict material properties for a wide-range of molecules. Schrödinger's AutoQSAR tools automates the generation of accurate QSAR models, which allows users to leverage machine learning tools without extensive background knowledge. For a complete description of how AutoQSAR automatically tests various models and makes selections, visit the Machine Learning for Materials Science tutorial.

DeepAutoQSAR integrates graph convolutional neural networks into the traditional AutoQSAR workflow, where DeepAutoQSAR treats a molecule as a graph consisting of nodes as atoms and edges as bonds. DeepAutoQSAR has been found to outperform traditional AutoQSAR for 'large' datasets (>5000 molecules) and perform similarly to traditional AutoQSAR for 'small' datasets (<5000 molecules) (see comparison here). A distinct advantage of DeepAutoQSAR is its ability to identify hidden patterns relevant to the property of interest through a series of convolution operations. You can read more about DeepAutoQSAR on our webpage as well as the references therein.

In this tutorial, we will use the DeepAutoQSAR panel in MS Maestro to create a machine learning model to predict rate constants for a radical reaction (reductive dehalogenation of aryl halide) catalyzed by a series of organometallic iridium complexes. The experimental data set is provided from a recent publication from Mdluhi *et al.* (High-throughput Synthesis and Screening of Iridium (III) Photocatalysts for the Fast and Chemoselective Dehalogenation of Aryl Bromides. DOI:10.1021/acscatal.0c02247). This experimental dataset explores a series of ~1000 $[Ir(C^\wedge N)_2(N^\wedge N)]^+$ photocatalysts (octahedral iridium complexes with three, bidentate ligands) and measures rate constants using high-throughput colorimetric monitoring.

Herein, we use a data set of 863 of the iridium complexes and the experimental rate constants to train and evaluate machine learning models. The DeepAutoQSAR panel is used to generate a model to predict rate constants by training on the structure of each Ir complex and the associated rate constant. To test the generalizability of the model, rate constants are predicted for an unseen set of 50 complexes. The overall workflow is summarized in *Figure 1*.
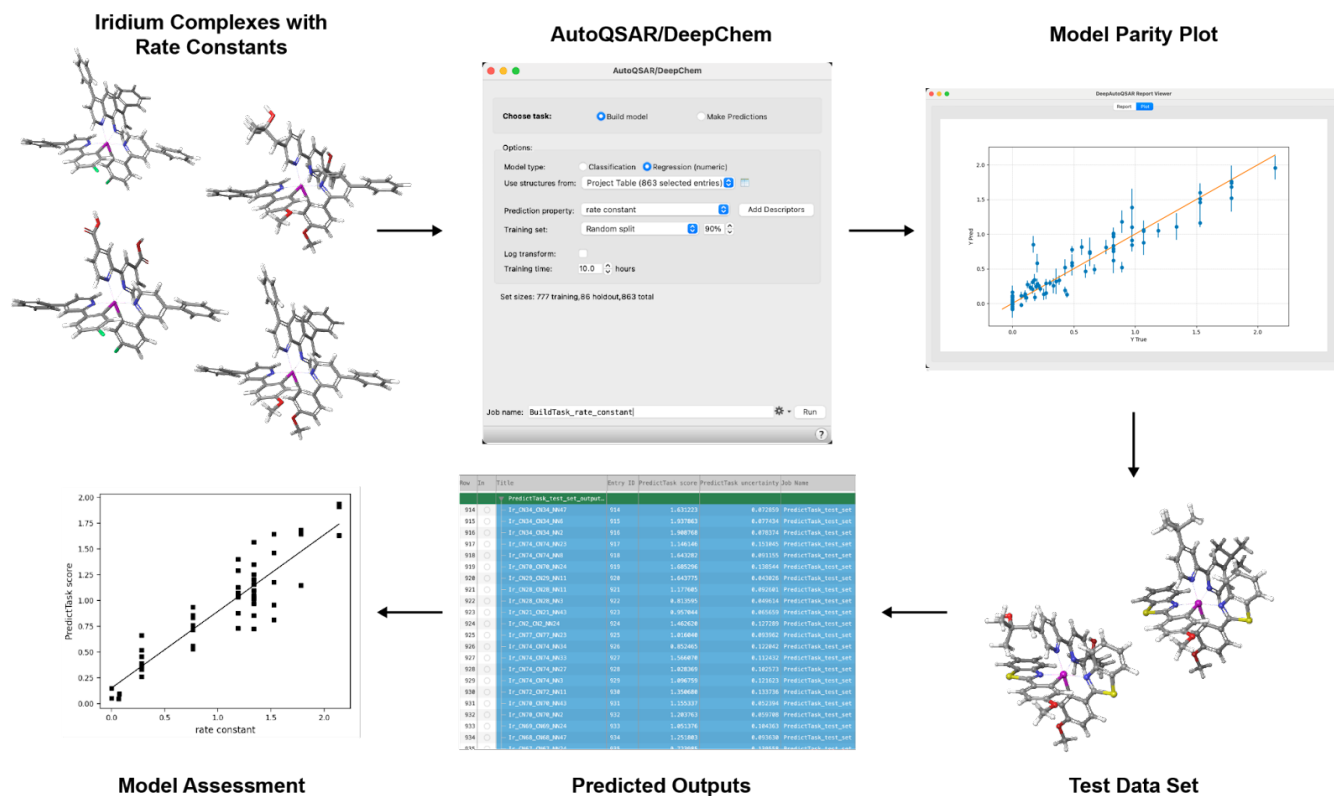
*Figure 1. Tutorial workflow showing the input Ir complexes, DeepAutoQSAR panel used to build machine learning models, and the output parity plot after model training. After training the model, an unseen test set was used to evaluate model performance. The workflow subsequently shows the Ir complexes, output predictions, and parity plot for the test set.*

Note that while this dataset is small enough that AutoQSAR could also be used, this tutorial focuses on using DeepAutoQSAR, which produces slightly more accurate predictions than traditional AutoQSAR.

For additional practice with the DeepAutoQSAR workflow, but with a categorical classification task, see the Machine Learning for Sweetness tutorial.

For additional practice with AutoQSAR, tutorials are available using the Materials Science Maestro suite to predict properties of small molecules, polymers and periodic systems: Machine Learning for Materials Science, Polymer Descriptors for Machine Learning and Periodic Descriptors for Inorganic Solids.
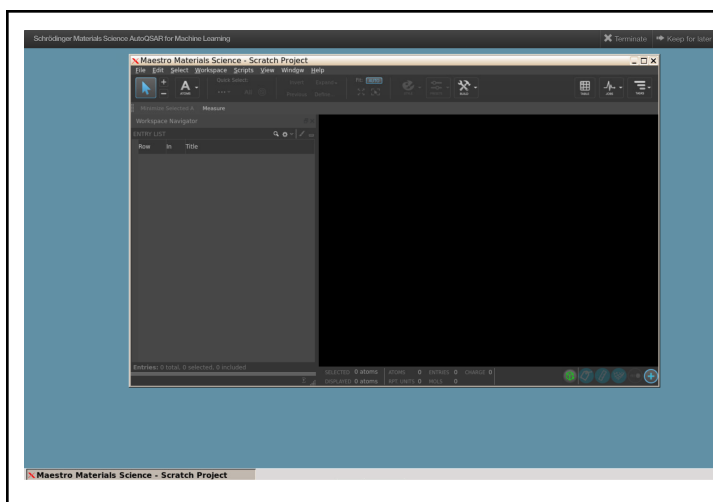
To learn about using pre-built machine learning models to predict volatility of organometallic complexes, please refer to the Machine Learning Property Prediction tutorial.

For alternative computational approaches for catalyst discovery, namely elucidating reaction mechanisms via various workflows, visit the Locating Transition States: Part 1 and Part 2 tutorials, as well as the Reaction Workflow for Polyethylene Insertion tutorial.

# 2. Creating Projects and Importing Structures

At the start of the session, change the file path to your chosen <u>Working Directory</u> in MS Maestro to make file navigation easier. Each session in MS Maestro begins with a default <u>Scratch Project</u>, which is not saved. A MS Maestro project stores all your data and has a `.prj` extension. A project may contain numerous entries corresponding to imported structures, as well as the output of modeling-related tasks. Once a project is saved, the project is automatically saved each time a change is made.

Structures can be built in MS Maestro or can be imported using **File > Import Structures** (or drag-and-dropped), and are added to the <u>Entry List</u> and <u>Project Table</u>. The <u>Entry List</u> is located to the left of the <u>Workspace</u>. The <u>Project Table</u> can be accessed by **Ctrl+T (Cmd+T)** or **Window > Project Table** if you would like to see an expanded view of your project data.

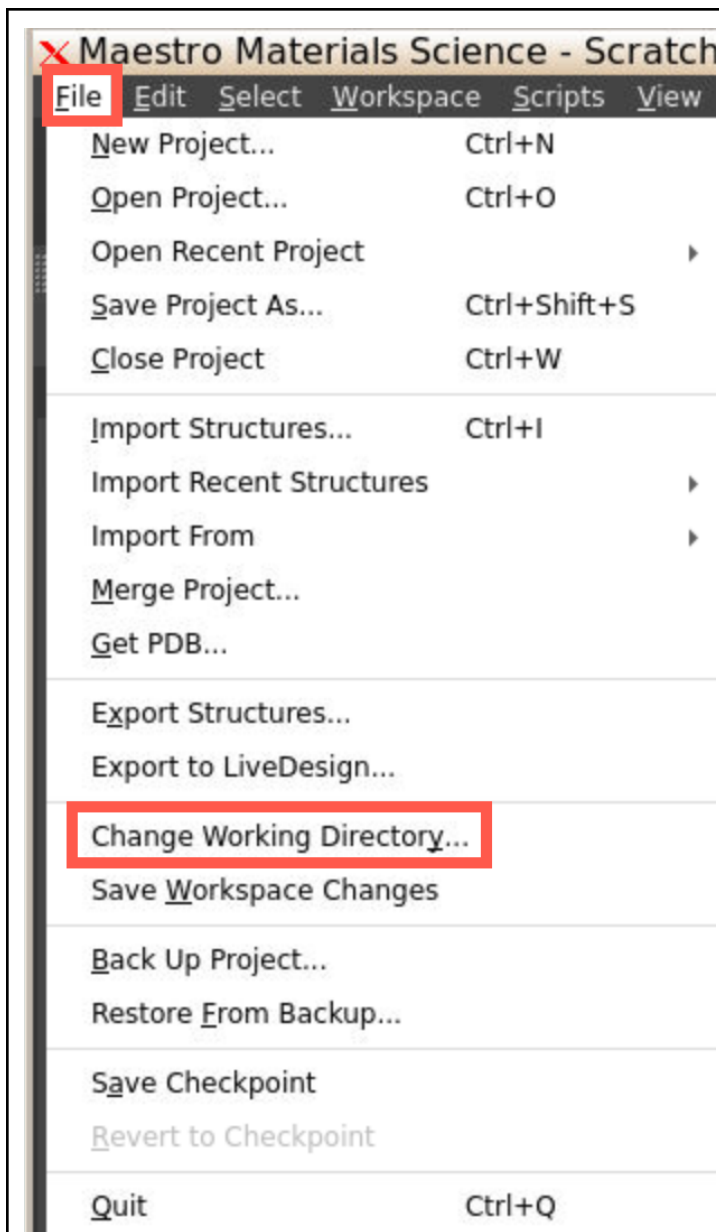| | |
|---|---|
|  | 1. Launch the tool in the nanoHUB interface<br>    ○ Launching the tool will automatically open up MS Maestro |

Figure 2-1. Change Working Directory option.
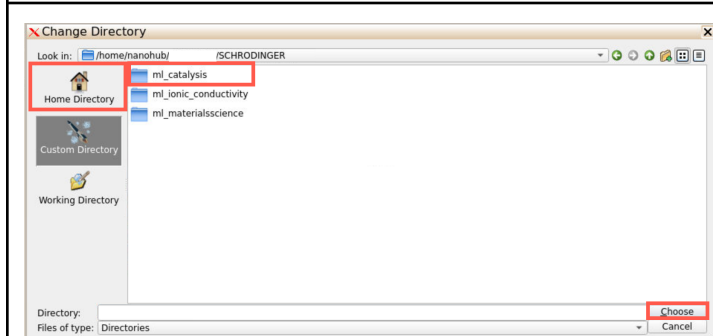
2. Go to **File > Change Working Directory**



Figure 2-2. Selecting the Working Directory.

1. Navigate to your Home Directory then the SCHRODINGER directory
2. Select *catalysis_ml*, and click **Choose**
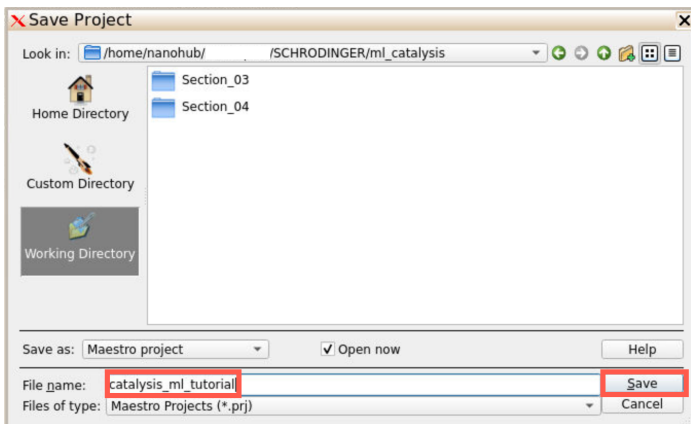   - All files needed to execute this tutorial are included in this directory

| | |
|---|---|
| <br>*Figure 2-3. Save Project panel.* | 3. Go to **File > Save Project As**<br>4. Change the *File name* to **catalysis_ml_tutorial**, click **Save**<br>    ○ The project is now named `catalysis_ml_tutorial.prj` |
| <br>*Figure 2-4. Import the starting structures.* | Let's import the data set:<br><br>5. Go to **File > Import Structures**<br>6. Navigate to where you downloaded the provided tutorial files (presumably in your <u>working directory</u>) and **choose** `train.mae` from the provided tutorial files<br>7. Click **Open** |
| <br>*Figure 2-5. The entry list and a complex after importing.* | The entry list is updated to include the 863 entries. Feel free to stylize and visualize any of the provided structures.<br><br>*Note:* The model complexes were prepared using Materials Science Maestro structure building capabilities (see the <u>Organometallic Complexes</u> tutorial for relevant workflows). |

*Figure 2-6. Viewing the rate constants in the Project Table.*

Each entry has a rate constant associated with it. These can be visualized in the Project Table (). Use the Property Tree () to add the rate constant property (under **All > Materials Science > Secondary > rate constant**)

# 3. Building a Machine Learning Model Using DeepAutoQSAR

In this section, we will use the DeepAutoQSAR panel to train a machine learning model for rate constant prediction. For a complete description of how AutoQSAR automatically tests various models and makes selections, visit the Machine Learning for Materials Science tutorial.
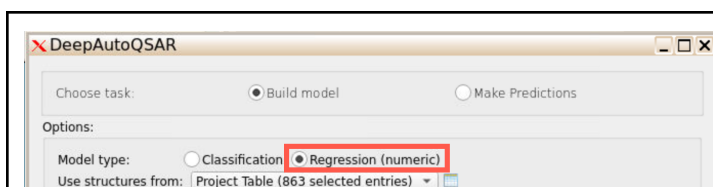


*Figure 3-1. Choosing task and options.*

1. Ensure that all 863 entries are selected from the entry list (use **Shift + Click** or click on the entry group header)
2. Go to **Tasks > Browse All > Discovery Informatics and QSAR > DeepAutoQSAR**
   - The DeepAutoQSAR panel opens
3. Ensure that *Build model* is checked
4. For *Model type*, choose **Regression**
   - Because the data is numerical and continuous, we use the regression *Model type*
5. Ensure that for *Use structures from*, **Project Table (863 selected entries)** is chosen
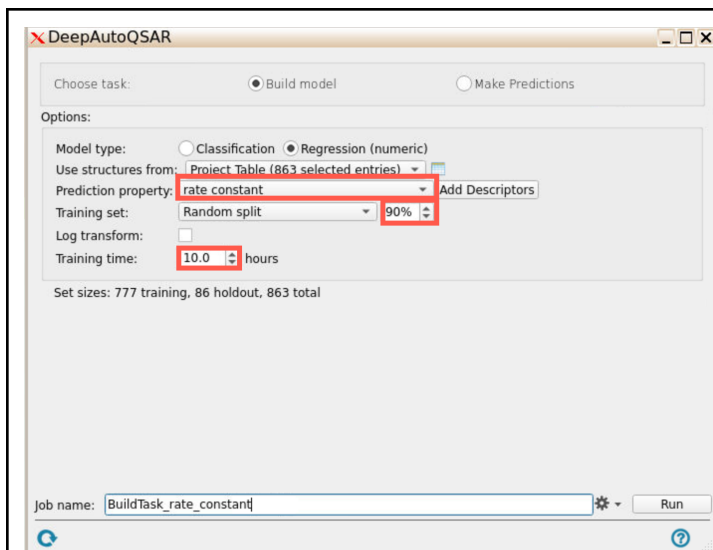
Figure 3-2. Defining the training.

6. Change the *Prediction property* dropdown to **rate constant**
7. Set **90%** for the *Random split*
   ○ This is the percentage of data to set aside between train and test sets, where 90% of the data is used to train the model and 10% of the data is used to test the model
   ○ This is a relatively large data set. The 90:10 split ensures that there is significantly more data in the training set than the test set, but still enough data in the test set to assess model performance
8. Set the *Training time* to **10 hours**
   ○ For datasets with >800 inputs like this one, a 10 hour training time is sufficient to ensure the best models are determined
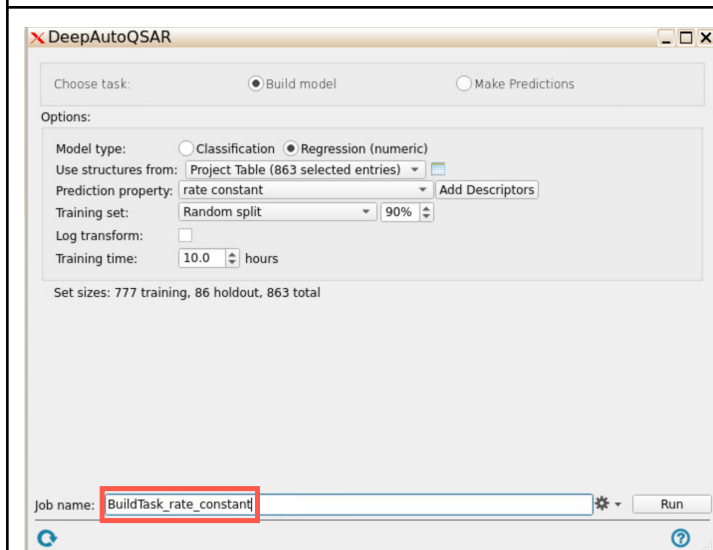


Figure 3-3. Naming the job.

9. Change the *Job name* to **BuildTask_rate_constant**

🛑 **DO NOT RUN**

The job would run for 10 hours as prescribed. The provided data is available for proceeding in **Section 4**. You can proceed to **Section 4** where steps are provided for importing the pre-computed models.

10. **Close** the DeepAutoQSAR panel (or simply move it to the side of your window – we will return to it in a moment)

# 4. Viewing the Machine Learning Model and Predicting

Using the DeepAutoQSAR panel, we can proceed to view the generated models, and use these to make predictions on an unseen data set.
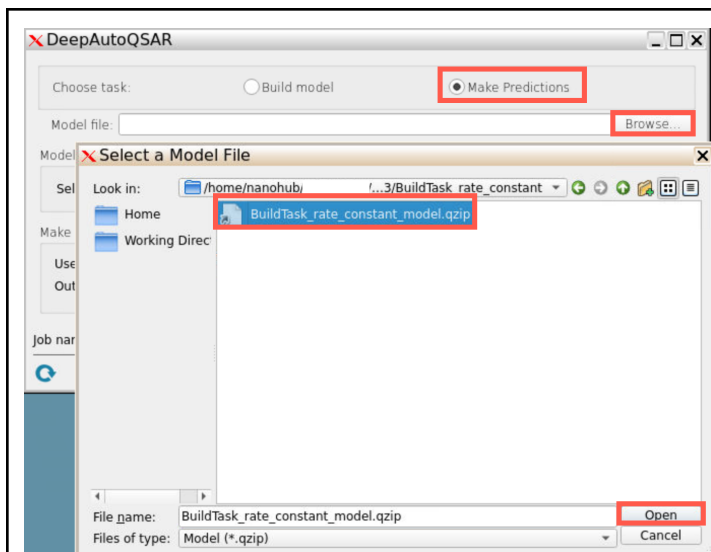
*Figure 4-1. Loading the .qzip file.*

The output can be analyzed and used for predictions back in the DeepAutoQSAR panel:

1. Return to **Tasks > Browse All > Discovery Informatics and QSAR > DeepAutoQSAR**
   - The DeepAutoQSAR panel opens
2. For *Choose task*, switch to **Make Predictions**
3. To choose the *Model file* click **Browse**, navigate to the `Section_03 > BuildTask_rate_constant > BuildTask_rate_constant_model.qzip` file and click **Open**
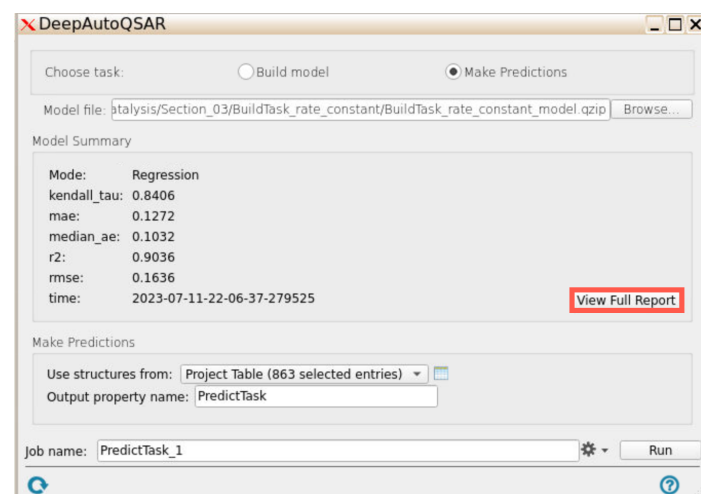   - The panel will parse the .qzip file and the *Model Summary* section will be populated



*Figure 4-2. Viewing the Model Summary.*

Begin by analyzing the *Model Summary* output. The data presented is a summary of the statistics of the model on the test set. We observe that the DeepAutoQSAR achieves a high $R^2$ of ~0.90 (denoted as r2) and low root-mean-squared error (rmse) of ~0.16 (an ideal model would have $R^2$ of 1 and RMSE of 0).

4. Click **View Full Report**



The *Report* tab includes a raw copy of the JSON output of DeepAutoQSAR. This report contains information on the top four best-performing model ensembles, including their metrics, the classification method used (e.g. dNN, random forest, etc.) and relevant model meta-parameters.

For a complete description of how AutoQSAR automatically tests various models and makes selections, visit the Machine Learning for
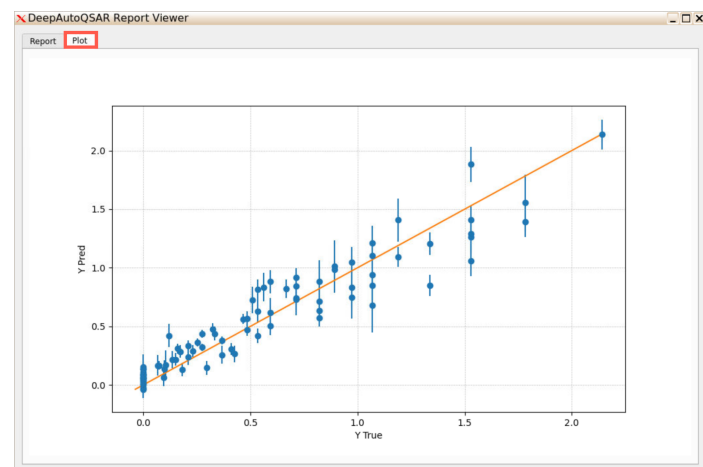
| | |
|---|---|
| *Figure 4-3. Viewing the Report tab.* | Materials Science tutorial. |
| <br><br>*Figure 4-4. Viewing the parity plot.* | 5. Click on the **Plot** tab<br><br>For regression models such as this, the *Plot* tab shows a parity plot.<br><br>6. **Close** the DeepAutoQSAR Report Viewer |
| <br><br>*Figure 4-5. Imported structures in the entry list.* | Now, we will use the trained model to make predictions on an unseen data set of iridium complexes that were not in the training data. These complexes have known rate constants from the same experimental study, which we can use to assess the quality of the model for making predictions outside the training set.<br><br>7. **Close** the DeepAutoQSAR panel (or simply move it to the side of your window – we will return to it in a moment)<br>8. Go to **File > Import Structures**<br>9. Navigate to where you downloaded the provided tutorial files (presumably in your working directory), choose `test.mae` and click **Open**<br>   ○ A new entry group is added to the entry list titled test (50) |
| <br><br>*Figure 4-6. The DeepAutoQSAR panel with the* | 10. Select the entire test (50) group from the entry list<br>   ○ Recall that select means to highlight the group in the entry list<br>11. Return to the DeepAutoQSAR panel<br>12. Ensure that the panel reflects the progress from the above steps: *Make Predictions* is selected, the .qzip file is loaded and the *Model Summary* is shown |

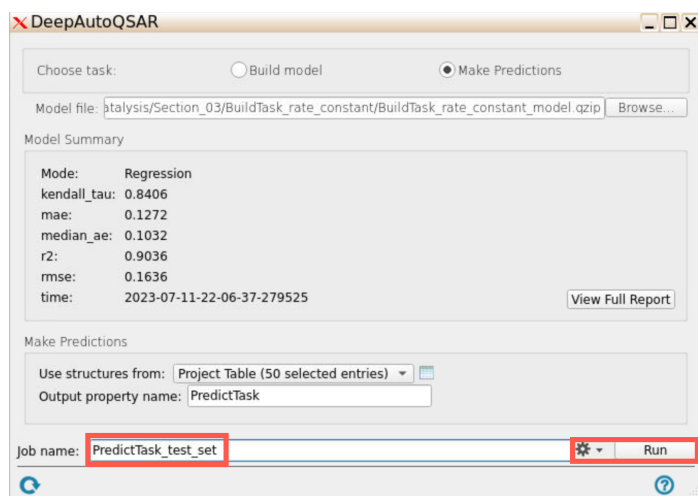| | |
|---|---|
| *prediction set selected.* | 13. In the *Make Predictions* section of the panel, ensure that **Project Table (50 selected entries)** is chosen for *Use structures from* |
| <br><br>*Figure 4-7. Naming and running the job.* | 14. For *Output property name*, maintain **PredictTask**<br>　○ This will be the name of the predicted property in the project table<br>15. Change the *Job name* to **PredictTask_test_set**<br><br>Adjust the job settings ( ⚙️▾ ) as needed. This job requires a Linux host. The job can be completed in about 5 minutes, which is of course many orders of magnitude faster than computing the rate constants of 50 systems from first principles. If you do not wish to run the job, feel free to simply import `Section_04 > PredictTask_test_set > PredictTask_test_set_output.maegz` from the provided tutorial files<br><br>16. Click **Run**<br>17. **Close** the DeepAutoQSAR panel |
| <br><br>*Figure 4-8. Viewing the output in the Project Table and opening the plots* | When the job is complete or after importing, a new entry group is added to the <u>entry list</u> titled PredictTask_test_set_output1 (50) containing the same 50 entries, but now with predicted rate constant property. The data can be analyzed in the <u>Project Table</u><br><br>18. Open the Project Table ( ⊞ TABLE )<br><br>19. Use the **Property Tree** ( ⊟ ) to include the *Predicted Task score* and *uncertainty* properties (**Check** the properties of interest under **All > Maestro > Predict Task score/uncertainty**) |

| | We can see predicted scores for the various molecules as well as uncertainty values.

To compare these values to the known values we will draw a scatter plot.

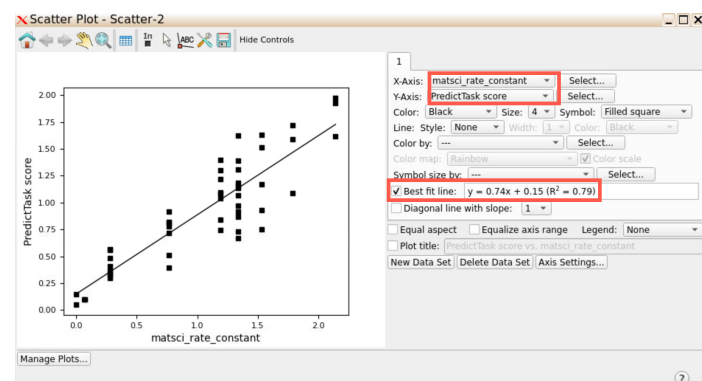20. Click the **Manage Plots** ( ) button |
|---|---|
| 

*Figure 4-9. A scatter plot of the predicted data versus the known values for the test set.* | 21. Click **New Scatter Plot**
22. For *X-Axis* select **matsci_rate_constant**
    ○ These are the actual values of the target property
23. For *Y-Axis* select **PredictTask score**
    ○ These are the ML predicted values
24. **Check** *Best fit line*
    ○ A regression line and equation is added

This scatter plot was generated from the values in `PredictTask_test_set_output.maegz`. If you performed this calculation your scatter plot and best fit line will differ slightly. |

This workflow highlights the computational efficiency achieved when using ML approaches as compared to other computational (*e.g. ab initio* calculations) or experimental approaches. While this tutorial uses a relatively small dataset, one could expect that a larger training set would further improve prediction accuracy.

# 5. Conclusion and References

In this tutorial, we learned how to use the DeepAutoQSAR panel to build machine learning models to predict experimentally determined rate constants for a series of iridium complexes. The DeepAutoQSAR model can generalize to unseen data sets and generate fast predictions (~seconds-minutes) as compared to *ab initio* or experimental measurements (~hours-days), enabling the screening of catalysts for enhanced reaction rates. While this tutorial focuses on reaction rate constants of iridium complexes, the workflow can be extended to other catalyst types and properties.

**For further learning:**

For introductory content, focused on navigating the Schrödinger Materials Science interface, an Introduction to Materials Science Maestro tutorial is available. Please visit

the [materials science training website](#) for access to 50+ tutorials. For scientific inquiries or technical troubleshooting, submit a ticket to our Technical Support Scientists at [help@schrodinger.com](mailto:help@schrodinger.com).

For self-paced, asynchronous, online courses in Materials Science modeling, including access to Schrödinger software, please visit the [Schrödinger Online Learning](#) portal on our website.

For some related practice, proceed to explore other relevant tutorials:

- For more machine learning:
    - [Machine Learning for Materials Science](#)
    - [Polymer Descriptors for Machine Learning](#)
    - [Periodic Descriptors for Inorganic Solids](#)
    - [Optoelectronics Active Learning](#)
    - [Machine Learning for Sweetness](#)
    - [Machine Learning Property Prediction](#)
    - [Machine Learning for Ionic Conductivity](#)
- For transition state searching with quantum mechanical methods in molecular or periodic systems:
    - [Locating Transition States: Part 1](#)
    - [Locating Transition States: Part 2](#)
    - [Reaction Workflow with Polyethylene Insertion](#)


**For further reading:**

- Help documentation on [DeepAutoQSAR](#)
- High-throughput Synthesis and Screening of Iridium(III) Photocatalysts for the Fast and Chemoselective Dehalogenation of Aryl Bromides. [DOI:10.1021/acscatal.0c02247](#)
- DeepAutoQSAR: Scalable, Intuitive, Deep-learning QSAR models for Big Data Applications (Schrödinger [white paper](#))
- DeepAutoQSAR Hardware Benchmark (Schrödinger [white paper](#))
- Design of Organic Electronic Materials With a Goal-Directed Generative Model Powered by Deep Neural Networks and High-Throughput Molecular Simulations. [DOI:10.3389/fchem.2021.800370](#)
- Active Learning Accelerates Design and Optimization of Hole-Transporting Materials for Organic Electronics. [DOI:10.3389/fchem.2021.800371](#)
- Some recent publications applying machine learning methods in catalysis and reactivity:
    - Machine Learning in Catalysis, From Proposal to Practicing. [DOI:10.1021/acsomega.9b03673](#)
    - Accelerated dinuclear palladium catalyst identification through unsupervised machine learning. [DOI:10.1126/science.abj0999](#)

- Univariate classification of phosphine ligation state and reactivity in cross-coupling catalysis. DOI:10.1126/science.abj4213
- Catalytic Performance of Cycloalkyl-Fused Aryliminopyridyl Nickel Complexes towards Ethylene Polymerization by QSPR Modeling. DOI:10.3390/catal11080920

# 6. Glossary of Terms

Entry List - a simplified view of the Project Table that allows you to perform basic operations such as selection and inclusion

Included - the entry is represented in the Workspace, the circle in the In column is blue

Project Table - displays the contents of a project and is also an interface for performing operations on selected entries, viewing properties, and organizing structures and data

Recent actions - This is a list of your recent actions, which you can use to reopen a panel, displayed below the Browse row. (Right-click to delete.)

Scratch Project - a temporary project in which work is not saved, closing a scratch project removes all current work and begins a new scratch project

Selected - (1) the atoms are chosen in the Workspace. These atoms are referred to as "the selection" or "the atom selection". Workspace operations are performed on the selected atoms. (2) The entry is chosen in the Entry List (and Project Table) and the row for the entry is highlighted. Project operations are performed on all selected entries

Working Directory - the location where files are saved

Workspace - the 3D display area in the center of the main window, where molecular structures are displayed