

SRSML24: STM Machine Learning Module

Steven R. Schofield

April 20, 2025

Overview

This module provides tools for machine learning analysis of scanning tunnelling microscopy (STM) data, including autoencoder models, clustering tools, and STM-specific preprocessing.

Getting the Code

Clone the repository from GitHub:

```
git clone https://github.com/srschofield/SRSML24.git
```

Installation

It is recommended to create a clean Python environment using `conda`. The following steps assume you are working on a macOS system with Apple Silicon:

```
# create and activate environment
conda create --name srsml24 python=3.8 -y
conda activate srsml24
```

```
# install packages
pip install -r requirements-macos.txt
```

Known Working Configuration

This module has been tested and is known to work with the following configuration on macOS 15.0.1 (Apple Silicon, M3 Pro chip):

- `python==3.8`
- `tensorflow-macos==2.16.2`
- `tensorflow-metal==1.1.0`
- `numpy==1.24.3`
- `pandas==1.5.3`
- `matplotlib==3.7.1`

- `scikit-learn==1.3.0`
- `scipy==1.10.1`
- `opencv-python==4.8.1.78`
- `Pillow==9.5.0`
- `joblib==1.3.2`
- `jupyter==1.0.0`
- `ipykernel==6.29.3`
- `keras-core==0.1.6`
- `spiepy==0.1.6`
- `access2thetmatrix==0.1.3`

These packages can be installed using the `requirements-macos.txt` file. The Python version is critical: other versions may cause compatibility issues with TensorFlow or other packages on Apple Silicon.

Python Files

- `data_prep.py` – Functions for data preparation, including slicing STM images into windows and saving them in efficient formats.
- `model.py` – Defines convolutional autoencoder and UNET-style models.
- `utils.py` – Utility functions for loading/saving models, feature arrays, and results.

Example Data and Scripts

- `example_data/` – Example STM data and latent feature arrays.
- `examples/` – Scripts demonstrating usage of the module, including training and inference workflows.

License

This work is licensed under the Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License (CC BY-NC-SA 4.0). You are free to:

- **Share** — copy and redistribute the material in any medium or format
- **Adapt** — remix, transform, and build upon the material

Under the following terms:

- **Attribution** — You must give appropriate credit, provide a link to the license, and indicate if changes were made.

- **NonCommercial** — You may not use the material for commercial purposes.
- **ShareAlike** — If you remix, transform, or build upon the material, you must distribute your contributions under the same license.

To view a copy of this license, visit: <https://creativecommons.org/licenses/by-nc-sa/4.0/>

Parameter Summary

Parameter	Description
General	
<code>job_name</code>	Label for the run, it will be the folder name for output.
<code>verbose</code>	If True , enables more detailed print output.
Matrix data file processing	
<code>flatten_method</code>	Method used to flatten STM images before analysis. Options are 'none', 'iterate_mask', 'poly_xy'.
<code>pixel_density</code>	All images will be converted to this pixel density (px/nm).
<code>pixel_ratio</code>	Images that have ratio of fast/slow scan direction less than this will be discarded. Setting to 1 means only complete (square) images are kept.
<code>data_scaling</code>	Multiplicative factor for z-height data. Setting to 1.e9 means that the range 0–1 (used for training) corresponds to 1 nm.
Window generation	
<code>window_size</code>	Side length of square image windows (in pixels).
<code>window_pitch</code>	Spacing between adjacent windows during tiling.
Data saving	
(Should remain defaults but options can be useful for examining data manually.)	
<code>save_windows</code>	If True , saves image windows as <code>.npy</code> files (True).
<code>together</code>	If True , saves windows per image in a single file (True).
<code>save_jpg</code>	If True , saves full STM images as JPGs (False).
<code>collate</code>	If True , flattens directory structure into one folder. (False).
<code>save_window_jpgs</code>	If True , saves image windows as JPGs. (False)
Autoencoder	
<code>model_name</code>	Label used to save and load the trained autoencoder model.
<code>batch_size</code>	Number of windows per training batch.
<code>buffer_size</code>	Size of shuffle buffer.
<code>learning_rate</code>	Learning rate for the optimizer.
<code>epochs</code>	Number of training epochs.
Clustering	
<code>cluster_model_name</code>	Name used when saving the clustering model.
<code>cluster_batch_size</code>	Number of latent vectors per clustering batch.
<code>cluster_buffer_size</code>	Size of buffer for clustering shuffle.
<code>num_clusters</code>	Number of clusters to form using KMeans.
<code>n_init</code>	Number of initializations for KMeans.
<code>max_iter</code>	Max iterations for KMeans convergence.
<code>reassignment_ratio</code>	Fraction of centroids reassigned each step.

Parameter	Description
Image prediction	
<code>predict_window_pitch</code>	Window spacing during prediction step.
<code>mtrx_train_data_limit</code>	Max number of training MTRX files to use.
<code>mtrx_test_data_limit</code>	Max number of validation MTRX files to use.
<code>train_data_limit</code>	Limit on number of training windows.
<code>test_data_limit</code>	Limit on number of validation windows.