

MetaMethyLib Usage Manual

Disclaimer: MetaMethyLib is currently available as an alpha release (work in progress). The current version is still under development.

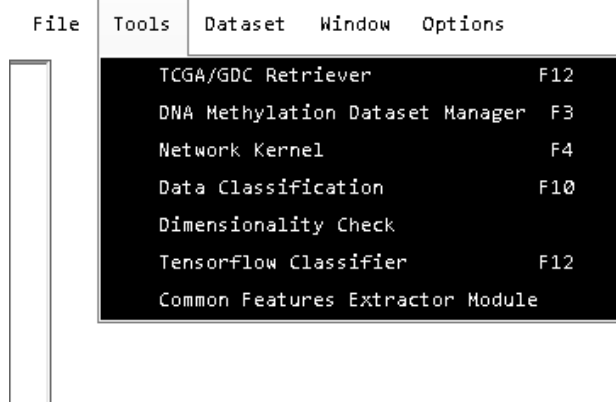
Main features:

- 1) Download DNA methylation data from TCGA portal.
- 2) Process the downloaded data to build the dataset.
- 3) Perform the metaheuristic feature selection of DNA methylation data (1st stage)
- 4) Perform the DNN pan-cancer classification (2nd stage).

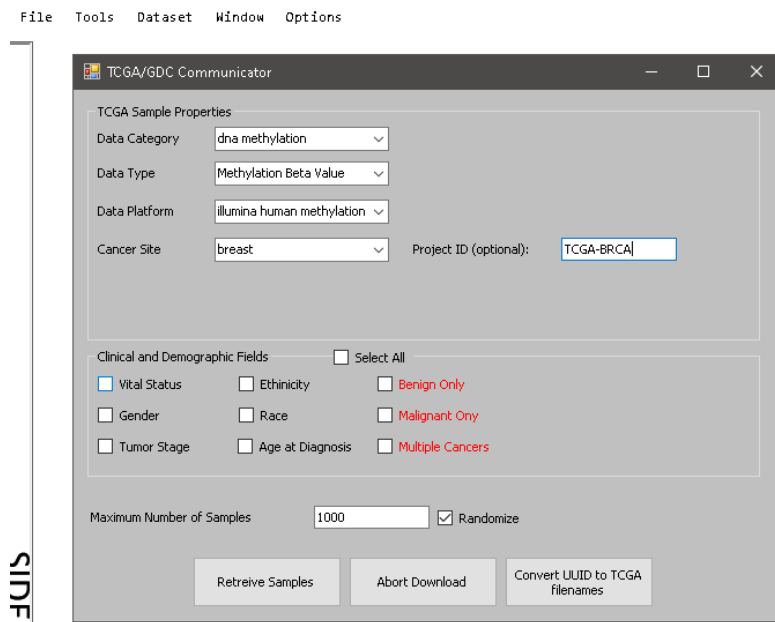
Data collection from TCGA portal

To download the DNA methylation data from the TCGA data portal, follow these steps:

- a. Launch the application and press F12 or navigate to **Tools->"TCGA/GDC Retriever"**



- b. Input the required information and then click on **"Retrieve Samples"** button.



Select a folder where the DNA methylation samples, along with the .csvs manifest file, will be saved.

Notes:

- 1) The optional project ID is to specify the exact TCGA project name if the current cancer site contains more than one project. You can select optional clinical field (e.g. vital status, gender, race, etc.) to be downloaded if available on the TCGA portal.
- 2) You can select the “Benign Only” option to download benign samples only, or “Malignant Only” to download malignant samples only. By default, the application will download random mixed samples (if available).
- 3) The “Convert UUID to TCGA filenames” is currently only required if the downloaded samples are based on the “Illumina 450k methylation platform”. You can use this to convert the filenames from UUID to TCGA ID, so that you can use the downloaded samples with the application.

Building the data set

To build the unified data set that will be used in the 1st stage, press F3 or navigate to **Tools->“DNA Methylation dataset manger”**. Then follow these steps:

The screenshot shows the 'MetaMethyLib Dataset Manager' window. It has a dark header bar with the title. Below the header, there are two tabs: 'Process Manifest File' and 'Process Samples'. The 'Process Manifest File' tab is active, showing a description: 'This module is used to process the methylation samples acquired from TCGA. All the samples should retain their original filenames and a manifest file describing all of them should be available.' Below this, there is a checkbox labeled 'Filter out Normal Tissues (Benign) ?' which is currently unchecked. The 'Process Samples' tab is also visible, showing a 'Dataset Creator' section with an 'Export Field' dropdown and a 'Choose the required metadata information to add with the beta value in the dataset' instruction. Below this, there is a 'Data Separator' dropdown set to 'Comma', a 'Sample-bValue Disposal Threshold (%)' slider set to 10%, a 'Convert nominal values to numerical equivalent ?' checkbox which is unchecked, and a 'Percentage of samples to process (%)' input field set to 100, with a 'starting at index' input field set to 0. There is an 'Export Dataset' button. At the bottom, there is a 'Samples Status' section with fields for 'Manifest File', 'Samples Folder', and 'Number of Processed Samples', all of which are currently empty. A 'Finish' button is located at the bottom right.

- a. Click on the **“Process Manifest File”** button and select the created CSV manifest file. Make sure you untick the **“Filter out Normal tissues”** checkbox if you need to process all entries (malignant + benign).
- b. Click on the **“Process Samples”** button and select the folder where the downloaded samples are located (illustrated in the figure below).

| | |
|--------------------------------------|--------------------------------------|
| 0a4dfb00-0321-4c63-9077-30de73d82fe3 | 0bfa648f-08ff-4935-b815-a14009d4ff02 |
| 3a31a682-a98c-4258-979b-659fb9ed0658 | 4b80bdda-d6df-4d1b-8815-f1445d244606 |
| 6fbb45ce-0260-4a81-9db8-0e855d8991c6 | 8aad1284-af8c-401a-81d4-10b68cb195cc |
| 9b8b68fa-5c3a-49c3-8eb7-d667f1b34d76 | 9d750eeb-5919-4275-9081-22c976714ef0 |
| 59a084f8-0af1-4607-a8af-566ca3cc4a7c | 59fc417e-cc24-481e-a76f-add8927a147b |
| 95efcdba-1c38-4501-9645-12d363a7872e | 131e715d-39be-40db-aa76-4390a9925378 |
| 965f5487-70fe-40e6-bcb0-315c567255e7 | 984e46be-661e-43d0-be5e-c8333ae36eea |
| 3497ba52-232e-4eb3-8a40-df254c03d60d | 6287eacc-c3c3-4e6f-9587-ecfb97cf9a6b |
| 111073ca-a5b0-418c-9499-ccb6931d0c3d | 270773a6-19a4-4e16-9a28-fc54fb1496cc |
| 2580713b-621a-44c3-bdba-c6d8daf5af9a | 5825414e-d0c2-4fbb-8bcd-fb15d472a420 |
| a4076529-4a7a-45bc-9ba4-985db75ae4a9 | aa70c164-e34f-4d30-aede-11c54bee4f8e |
| bd2e352f-c951-4bd6-bf72-cb195acfec87 | c48f6db1-fd5d-41ea-9e6a-40da6a25cab6 |
| d117d1ed-98be-401e-9bdb-26ce3c7b5fef | d3132c44-ba06-4090-b026-91b51a09b236 |
| e895d07a-19f7-449f-a145-d24b907944f6 | e4321f45-7ec1-48c8-8d1a-9efe8c9438fe |

- c. The application will start processing the data

```

MetaMethylLib Console
11.37% Processed sample: 0/73:01204c62-a786-4eec-a249-9f6a9f3a7c33
12.74% Processed sample: 37/73:8aad1284-af8c-401a-81d4-10b68cb195cc
14.11% Processed sample: 1/73:03f170ea-3243-4f3a-89a6-47a8013b6c73
15.48% Processed sample: 38/73:8b88a52d-5899-40c9-91de-764c70c131e6
16.85% Processed sample: 2/73:0a4dfb00-0321-4c63-9077-30de73d82fe3
18.22% Processed sample: 39/73:8f43cd26-16d2-47fc-a9f1-4f6e28e2df9b
19.59% Processed sample: 3/73:0bfa648f-08ff-4935-b815-a14009d4ff02
20.96% Processed sample: 40/73:930d7cef-1456-4d1c-b777-338be9297935
22.33% Processed sample: 4/73:111073ca-a5b0-418c-9499-ccb6931d0c3d
23.70% Processed sample: 41/73:95efcdba-1c38-4501-9645-12d363a7872e
25.07% Processed sample: 5/73:131e715d-39be-40db-aa76-4390a9925378
26.44% Processed sample: 42/73:965f5487-70fe-40e6-bcb0-315c567255e7
27.81% Processed sample: 6/73:16303e25-30f1-41f8-83d7-834f2927aa5b
29.18% Processed sample: 43/73:984e46be-661e-43d0-be5e-c8333ae36eea
30.55% Processed sample: 7/73:1ea50b76-ec5e-4cee-9ef1-953e2bhc7c96
31.92% Processed sample: 44/73:9a93c794-7cf8-43b6-bbc6-89dd29263d36
33.29% Processed sample: 8/73:2461bbab-9c64-4464-a0f1-aeb629d6c1c1
34.66% Processed sample: 45/73:9b8b68fa-5c3a-49c3-8eb7-d667f1b34d76
36.03% Processed sample: 9/73:2479df72-341e-4488-89fa-444ffaa7fc74
37.40% Processed sample: 46/73:9d750eeb-5919-4275-9081-22c976714ef0
38.77% Processed sample: 10/73:248fc4e3-cd38-4759-be15-975e45ed60f5
40.14% Processed sample: 47/73:a1873c8b-73f8-41e0-a79e-057cde49057
41.51% Processed sample: 11/73:2580713b-621a-44c3-bdba-c6d8daf5af9a
42.88% Processed sample: 48/73:a234b44f-e7ec-45a8-8436-f49eb062c3c5
44.25% Processed sample: 12/73:270773a6-19a4-4e16-9a28-fc54fb1496cc
45.62% Processed sample: 49/73:a4076529-4a7a-45bc-9ba4-985db75ae4a9
46.99% Processed sample: 13/73:2c89417a-fa80-4ef6-9a5b-0178a9hd0ff9
48.36% Processed sample: 50/73:a65fdb7b-35b6-4a87-8ad1-337b4a941856
49.73% Processed sample: 14/73:3497ba52-232e-4eb3-8a40-df254c03d60d
  
```

- d. After the data is successfully processed, the list of available clinical information will be displayed (if available). Select any clinical information that you want to include in the data set file

MetaMethylLib Dataset Manager

Process Manifest File This module is used to process the methylation samples acquired from TCGA. All the samples should retain their original filenames and a manifest file describing all of them should be available.

Process Samples ☐ Filter out Normal Tissues (Benign)?

Dataset Creator
Export Field Choose the required metadata information to add with the beta value in the dataset

☐ analysis.workflow_type ☐ case_id

☐ ethnicity ☐ gender

☐ race ☐ vital_status

☐ diagnoses.0.age_at_diagnosis ☐ diagnoses.0.tumor_stage

Data Separator Comma

Sample-bValue Disposal Threshold (%) 10%

Convert nominal values to numerical equivalent? ☐

Percentage of samples to process (%) 100 starting at index 0

Export Dataset

Samples Status

Manifest File D:\PhD\Datasets\Original Cancer Datasets (Benign+Tumor)\STOMACH\Manifest.csv

Samples Folder D:\PhD\Datasets\Original Cancer Datasets (Benign+Tumor)\STOMACH

Number of Processed Samples 73

Finish

- e. Select the cut-off threshold of missing beta values from the original TCGA DNA methylation samples. All the samples that contain missing beta values higher than the specified threshold will not be included in the final data set. You can also specify the percentage of samples that you need to include in the final data set (e.g. 70%, 30%).
- f. The application will create the data set file (.mml) that will be used later, in the metaheuristic feature selection stage.

Implementing the 1st Stage (Metaheuristic Feature Selection)

- a. First, you need to load the generated data set file (.mml). You can press “Ctrl+D”, or navigate to **Dataset->”Load dataset file (.mml)”**



- b. Then, you need to start a new study. To do so, press F5 or navigate to **File->New->Study**. Enter the desired study name and click Finish.

The 'MetaMethyLib Input' dialog box is shown. It has a title bar 'MetaMethyLib Input' and a subtitle 'Enter a name for this study'. Below the subtitle is a text input field containing 'Stomach_Study_1'. At the bottom of the dialog is a 'Finish' button.

- c. Enter the data for the nested genetic algorithms(s) and click on the “**Run Study**” button. (note: the animation is just to give confidence that the convergence is not locked in a non-responsive state).

The screenshot shows the 'Stomach_Study_1' application window. The title bar is 'Stomach_Study_1'. The window has a menu bar with 'File', 'Tools', 'Dataset', 'Window', and 'Options'. The main area is divided into two sections: 'Genetic Algorithm Parameters' and 'Console'. The 'Genetic Algorithm Parameters' section has two sub-sections: 'Outer Genetic Algorithm' and 'Inner Genetic Algorithm'. The 'Outer Genetic Algorithm' section has input fields for 'Crossover Probability (%)' (30), 'Mutation Probability (%)' (5), 'Population Size' (1000), and 'Max Generations' (100). The 'Inner Genetic Algorithm' section has input fields for 'Crossover Probability (%)' (30), 'Mutation Probability (%)' (5), 'Population Size' (1000), and 'Max Generations' (200). The 'Console' section shows a log of the application's execution, including the processing of TCGA samples and the start of the genetic algorithm. The 'Run Study' button is highlighted. To the left of the window, the text 'SIDE PANEL' is visible. At the bottom right, there is a 'Stomach_Study_1 Generation Tracker, Generation 0' table with columns for 'Chromosome 0', 'Chromosome 1', and 'Chromosome 2'.

| Chromosome 0 | Chromosome 1 | Chromosome 2 |
|--------------|--------------|--------------|
| | | |
| | | |
| | | |

- d. The first stage will start its convergence process as specified. It is important to note that you should “save” the study first before running it, to make sure that the auto-save feature will work. This feature exports the current progress to the workspace file after each generation. This file can be loaded anytime, and the application can resume the convergence process from where it stopped. In order to save the workspace (including all the active study sheets), navigate to **File->”Save workspace”** or process “Ctrl+S”.
- e. The following files will be created and updated throughout the convergence process.

| Name | Date modified |
|--|----------------|
| InnerChromosomeClustersContents.mmli3c | 5/17/2021 3:00 |
| InnerChromosomeClustersNumber.mmlicn | 5/17/2021 3:00 |
| InnerChromosomeClustersSizes.mmliccs | 5/17/2021 3:00 |
| OuterChromosomeContent.mmlocc | 5/17/2021 3:00 |
| OuterChromosomeFitness.mmlocf | 5/17/2021 3:00 |

Notes: by default, the application only uses local convergence on the host PC. If you need to use the network convergence feature. You need to launch the application in daemon mode and preferably, pass the required data set. To launch the application on network daemon mode, you need to run it from a command prompt (CMD) or by creating a shortcut and adding the following argument

-ParallelDaemon

The network daemon will start listening on the default IPv4 address of the local host. You can specify the listening port and the data set file by adding the following arguments, respectively:

PORT:xxxx MMLD:yyyy

Where ‘xxxx’ is the port number and ‘yyyy’ is the full path of the .mmlD data set file.

For example, to run the application as network daemon on port 20000 and data set file located in “d:\project\stomach.mmlD”, you should use the following command:

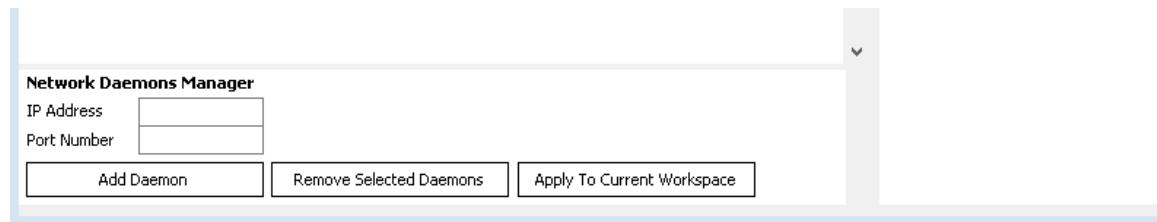
MetaMethyLib.exe -ParallelDaemon PORT:20000 MMLD:d:\project\stomach.mmlD

```

Select MetaMethyLib Console
MetaMethyLib.exe -ParallelDaemon PORT:20000 MMLD:d:\project\stomach.mmlD
<Running in Network Daemon Mode>>
dataset file d:\project\stomach.mmlD exists !
MetaMethyLib Notification
Loading MMLD Dataset, Check console for debug messages...
processing TCGA sample TCGA-BR-4367-11A
processing TCGA sample TCGA-BR-4183-01A

```

To connect the host application to the running daemon(s), navigate to **Tools->"Network Kernel"**, or press F4.



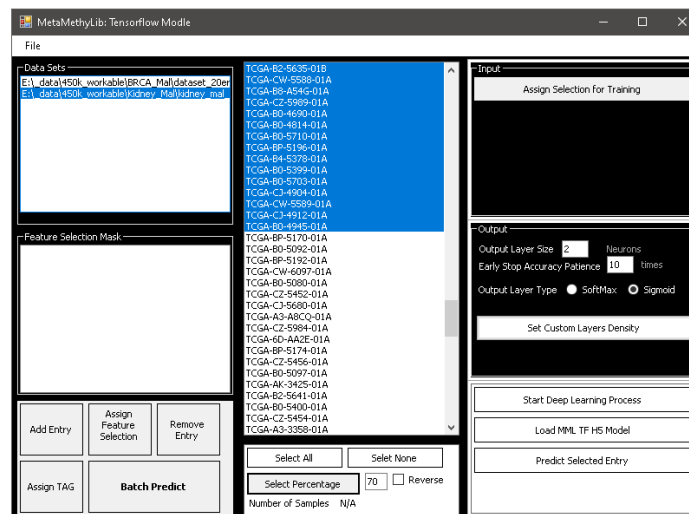
Write the IP address and port number of the daemon(s) and click the **"Add Daemon"** button and then click on the **"Apply to Current Workspace"** button.

Implementing the 2nd Stage (Classification)

For this stage, you need to make sure that all required pre-requisites are installed:

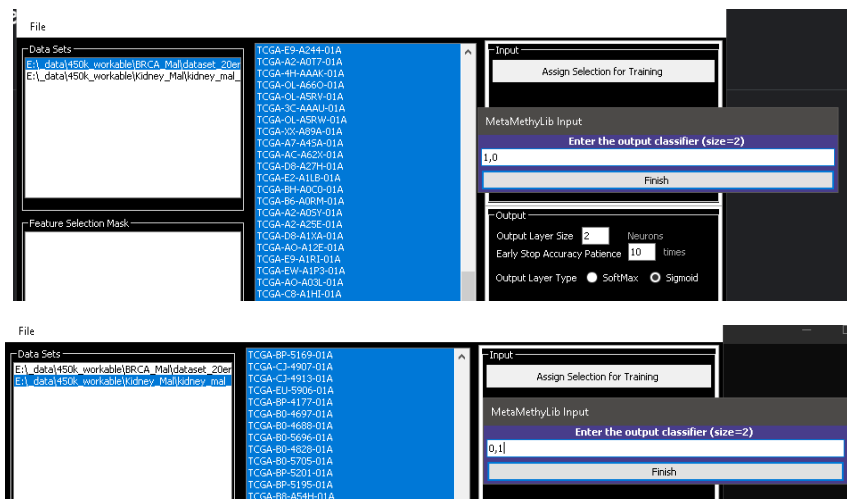
1. **Microsoft Visual C++ 2015-2019 Redistributable x86-x64** (download from <https://docs.microsoft.com/en-us/cpp/windows/latest-supported-vc-redist?view=msvc-170>)
2. **Python 3.8** (download from <https://www.python.org/ftp/python/3.8.0/python-3.8.0-amd64.exe>)
3. **Numpy** for python (pip) (install by running **'pip install numpy'** from any command prompt)
4. **Tensorflow** for python (pip) (install by running **'pip install tensorflow'** from any command prompt)
5. **Nvidia CUDA 11.2 x64** (download from <https://developer.nvidia.com/cuda-11.2.0-download-archive>)
6. **cuDNN 11.2 x64** (download from https://developer.nvidia.com/compute/machine-learning/cudnn/secure/8.1.0.77/11.2_20210127/cudnn-11.2-windows-x64-v8.1.0.77.zip)

To start the classification process, navigate to **Tools->"Tensorflow Classifier"** or press F12.



- a. First, add the required data set for each cancer type. You can do so by clicking the **"Add entry"** button. The data set(s) are the same ones generated during the 1st stage (.mmld file).

- b. Manually select the samples that you need to include in the training of the DNN, or alternatively, enter the required percentage (e.g. 70%, 30%, 20% etc.) in the percentage textbox and then click the **“Select Percentage”** button.
- c. If you need to assign different names to the training labels (instead of using the filename for each added cancer type), you can click the **“Assign TAG”** button for each selected entry.
- d. At this point, you can assign the generated features from the 1st stage to the entries by clicking on the **“Assign Feature Selection”** button. The extracted features from the 1st stage are found in the “.mmlocc” file. (To extract the common features from all the studies, you can use the common features extractor module by navigating to **Tools->“Common Features Extractor Module”**).
- e. Now, to define the output layer’s classes, select each entry, and click the **“Assign Selection for training”** button. Enter the One hot encoded data, separated by commas. You need to repeat this process for each one of the selected entries.



By default, the output layer size is estimated based on the added entries. You can explicitly change this number by modifying the value in **“Output Layer Size”** textbox. (e.g. binary classification).

- f. You can set an early termination threshold to stop the model from over fitting the data. You can also specify whether the output layer is based on a “sigmoid” or a “softmax” activation function.
- g. Finally, click on the **“Start Deep Learning Process”** button to start training the model using Keras Tensorflow.
- h. After the model has been successfully trained, the application will prompt you to save it on your local storage device. You can cancel this step if you do not want to save the trained model.
- i. To predict the result of any sample, first click on the entry in the “data sets” itembox and then select the desired sample by clicking on it. Afterwards, click on the **“Predict Selected Entry”** button.
- j. For convenience, you can also batch predict all the samples in the selected entry, by first selecting the desired entry, followed by clicking on the **“Batch Predict”** button. In this

case, after the batch prediction is complete, you will be prompted to save the results on your local storage device.