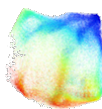


“How To EmbedSOM”



EmbedSOM plugin for FlowJo

EmbedSOM¹ is a dimensionality reduction algorithm that helps to quickly visualize your high-dimensional data in two dimensions. The design is based off FlowSOM², which is another popular unsupervised manifold learning algorithm that uses self-organizing maps (SOMs) to find clusters in the dataset. EmbedSOM uses the same manifold learning method to gain information about an approximate manifold that describes the high-dimensional space, and uses it to quickly compute a low-dimensional embedding to help visualize the dataset.

The v2.0 update adds new options for the EmbedSOM algorithm in addition to its previous functions, which allow users to load in a previously computed SOM from the FlowSOM clustering plugin. EmbedSOM will now compute its own embedding to generate a SOM as well as allow for the use of tSNE or UMAP landmark generation which can be used in place of a SOM to quickly generate the low-dimensional visualization.

New to version 2.1 is the ability to perform a semi supervised or supervised UMAP embedding using clustering labels from a previously calculated clustering algorithm. If enabled, the UMAP³ algorithm will use the clustering labels to help aid in learning and classification of the dataset to generate an embedding that shows generally well separated islands of cells. When selecting a population to perform a supervised embedding on which only a subset of events were clustered, not all events will have cluster labels. This is considered a semi supervised embedding and still generates nice results as the labeled data aids in classification while UMAP still captures structure in the unlabeled data.

The EmbedSOM algorithm has been developed and implemented as a R package by Miroslav Kratochvíl & Petr Ryšavý and the source code is available here: <https://github.com/xaexa/EmbedSOM>.

The R implementation of Uniform Manifold Approximation and Projection (UMAP) is available in the [uwot R package](#).

Download and installation

1. Place the plugin .jar file in your Plugins folder, and direct FlowJo to that folder using the Diagnostics section of the Preferences.
2. Make sure you have R installed and the R path is specified in the R Path field of the Diagnostics section of the Preferences.

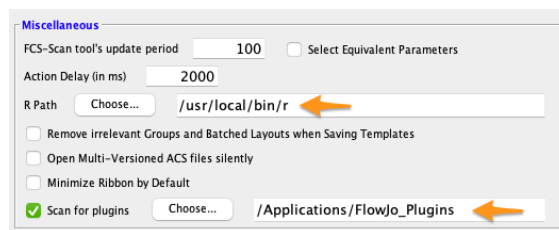


Figure 1 - The R path shown above is for MacOS. A typical Windows R path is C:\Program Files\R\R-4.2.2\bin\x64 for R version 4.2.2. A path to the plugins folder can also be specified.

- Running the plugin for the first time will install the needed R packages to allow the calculation to run in the R environment. Sometimes these installations can fail and you may need to manually install the R packages. To install the required R packages, use the following commands in R:

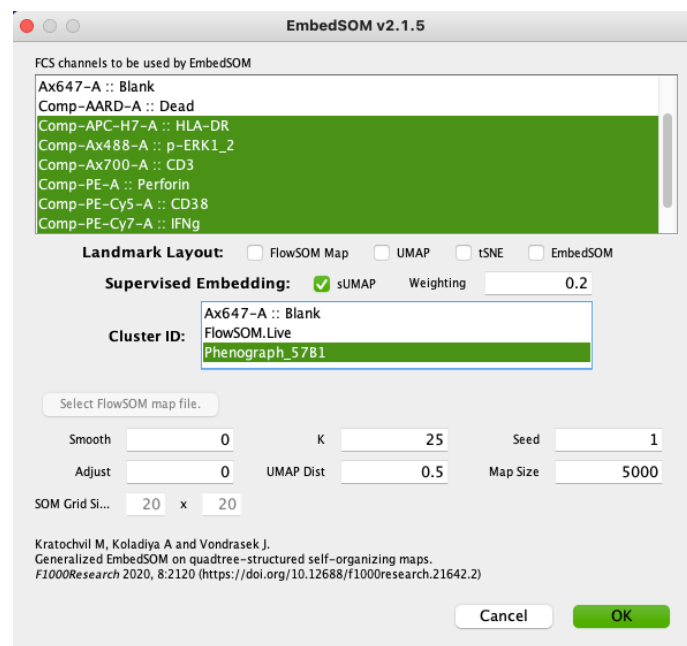
```
install.packages(c("ggplot2", "FNN", "igraph", "Matrix", "Rtsne", "umap",
                  "uwot", "utils", "devtools", "dplyr"))

devtools::install_github('exaexa/EmbedSOM')
```

Note This plugin was tested in R versions 4.1 & 4.2 and EmbedSOM R package version 2.1.2.

Usage

To run the EmbedSOM plugin on your FCS file, select the population of interest within the workspace. Go to the Workspace tab and select the EmbedSOM option from within the Plugins drop-down menu. This will bring up a dialog where researchers can choose which compensated parameters to use for dimensionality reduction. There are a variety of other options available to customize the embedding. Mousing over the fields will bring up a tool tip which describes the field's function.



Landmark Layout options in EmbedSOM Pick from different landmark generating functions to visualize the data in different ways. If you have previously used the FlowSOM plugin to calculate clustering results, you can use the SOM from FlowSOM to embed the same data or even map a new dataset to the SOM. Select the **FlowSOM Map** option then use the **Select FlowSOM map file** button to load in the FlowSOM .RData file that contains the SOM landmarks. This option will allow you to align multiple population positions across different files, as long as they have the exact same parameters. For example, we can concatenate multiple files together in FlowJo and run FlowSOM clustering, then use the FlowSOM .RData file to embed the same data into two dimensions. The same FlowSOM landmarks can be used again to embed new datasets to the same FlowSOM object and align these populations.

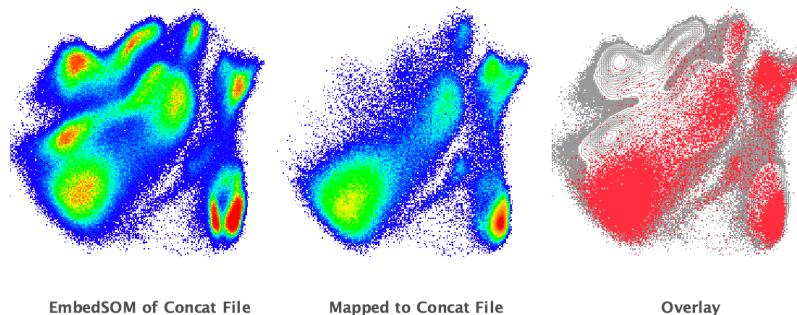


Figure 2 - Embedding on concatenated file (left-most image) consisting of different stim conditions is generated with FlowSOM .RData file and SOM landmarks. The middle image shows a new dataset mapped to the same FlowSOM landmarks and the results are overlaid in the right-most image.

Selecting between the **UMAP** and **tSNE** option will allow you to use these advanced dimensionality reduction methods to quickly generate the landmark positions using a subset of random events. The number of random input points to use is entered in the **Map Size** input field in the plugin UI. **K** is the number of nearest points used for the approximation, the default is set to 25 but can be raised up to 100 when embedding larger datasets. When using UMAP to generate landmarks the **UMAP Dist** is used to specify the minimum distance between points in low-dimensional space. Smaller values will push points together, while larger values will push points further apart in the final embedding.

For any of the Landmark options selected, increase the **Smooth** parameter to produce “smoother” but possibly more convoluted embedding. Default value 0 is adjusted to be good for most datasets; lower values (-1, -2,...) produce sharper embeddings with less focus on projection. Higher values (1,2,...) produce smoother embeddings that better explains the large-scale structure of data, but some small-scale details may get smoothed out. The **Adjust** parameter is a negative power factor for reducing the effect of non-local relevance measure on the outcome. Use 0 for plain projection; values above 100 usually push cells closer to respective SOM vertex positions.

Selecting the **EmbedSOM** option will use this improved embedding enrichment method to produce a smooth low-dimensional embedding using the grid size specified in the **SOM Grid Size** input fields. The **Seed** value can be used to create reproducible results by using the same seed value for the selected dataset.

Below we see embeddings created with each of the available methods in this plugin. The left-most plot was created by EmbedSOM SOM manifolds and shows a more smoothed-out embedding but still separates out the major lineages as shown by the colored islands. The next plot to the right was created using the FlowSOM .RData file and used the SOM landmarks created from the FlowSOM plugin and as map and embedded the new data. This embedding separates nicely and with the .RData file, can be applied to other sample files from the same experiment with the same parameter names. The UMAP result shows tighter grouping of islands, depicting more of the global structure of the data. The embedding on the right which used tSNE landmarks shows more local structure in the data but still separates the islands well thanks to the tunable parameters in this plugin.

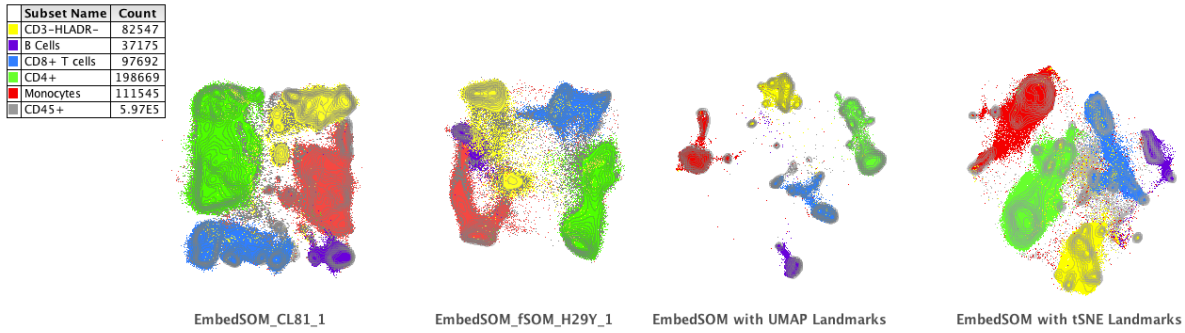


Figure 3 - Comparisons of the different embedding options on a CD45+ population is overlaid with major lineage populations. Each option is able to quickly produce embeddings to help visualize the data.

Supervised UMAP

Supervised Embedding options in EmbedSOM. When selecting the **sUMAP** checkbox, the relevant options will become actionable in the plugin ui. You have control over the weighting of your labels and which derived parameter to use as labels and drive the supervision.

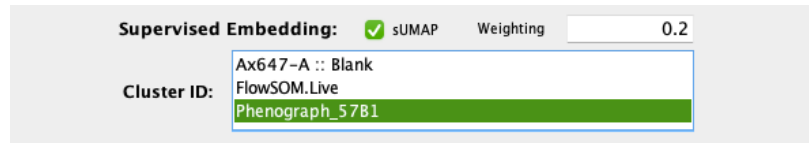


Figure 4 - The Supervised UMAP section of the plugin ui. You must select a clustering id or derived parameter here.

Enter a value in the **Weighting** field or leave the default setting of 0.2. This value is the relative weight of the supervised parameter. A higher value will use more supervision and drive the events into tighter groups. A lower value will deemphasize the effect of cluster input on the resulting low dimensional embedding. Please note that setting this value too high can sometimes have the opposite effect and the groups become spread out and do not represent the labels well. In most cases, the default value is sufficient to drive supervision and separate the classes of cells. Next, select a **Cluster ID** parameter from a previous clustering analysis in FlowJo. If you do not select the Cluster ID before clicking OK, you will be prompted again to select this derived parameter before you can continue.

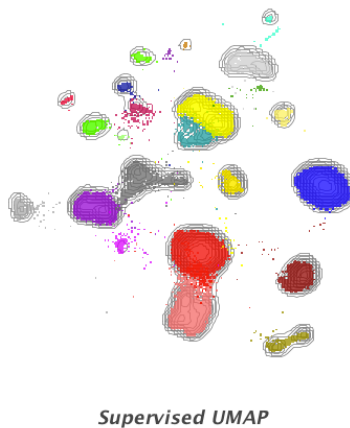


Figure 5 - UMAP embedding using Phenograph clusters as labels. The embedding of the same CD45+ population as shown in **Figure 3** now shows cells being classified into tighter groups in the final embedding.

Leave us your feedback

Please write to flowjo@bd.com with any questions or concerns.

References

1. Miroslav Kratochvíl, Abhishek Koladiya, and Jiří Vondrášek. “Generalized EmbedSOM on quadtree-structured self-organizing maps” F1000Research 8 (2019). doi:[10.12688/f1000research.21642.2](https://doi.org/10.12688/f1000research.21642.2)
2. Van Gassen S, Callebaut B, Van Helden MJ, Lambrecht BN, Demeester P, Dhane T, Saeys Y. FlowSOM: Using self-organizing maps for visualization and interpretation of cytometry data. *Cytometry A*. 2015; [87\(7\):636-45](#).
3. Examples of using supervised and semisupervised UMAP can be found [here](#).