

Sofie Van Gassen, Britt Callebaut, Tom Dhaene and Yvan Saeys

Ghent University

March, 2017

## Abstract

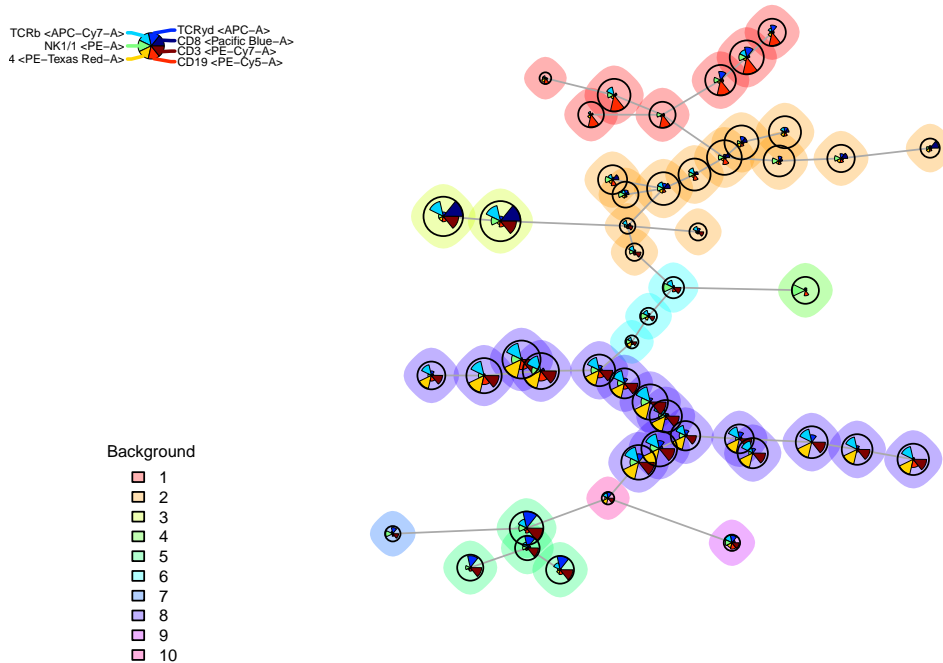
The *FlowSOM* package provides new visualization opportunities for cytometry data. A four-step algorithm is provided: first, the data is read and preprocessed, then a self-organizing map is trained and a minimal spanning tree is built, and finally, a meta-clustering is computed. Several plotting options are available, using star charts to visualize marker intensities and pie charts to visualize correspondence with manual gating results or other automatic clustering results.

## 1 The easy way

---

The easiest way to use this package is using the wrapper function `FlowSOM`. It has less options than using the separate functions, but in general it has enough power. It returns a list, of which the first item is the FlowSOM object (as required as input by many functions in this package) and the second item is the result of the metaclustering.

```
> library(FlowSOM)
> fileName <- system.file("extdata", "lymphocytes.fcs",
+                          package="FlowSOM")
> fSOM <- FlowSOM(fileName,
+                 # Input options:
+                 compensate = TRUE, transform = TRUE, toTransform=c(8:18),
+                 scale = TRUE,
+                 # SOM options:
+                 colsToUse = c(9,12,14:18), xdim = 7, ydim = 7,
+                 # Metaclustering options:
+                 nClus = 10,
+                 # Seed for reproducible results:
+                 seed = 42)
> PlotStars(fSOM$FlowSOM,
+           backgroundValues = as.factor(fSOM$metaclustering))
```



## 2 Reading the data

The FlowSOM package has several input options.

The first possibility is to use an array of character strings, specifying paths to files or directories. When given a path to a directory, all files in the directory will be considered. This process does not happen recursively. You can specify a pattern to use only a selection of the files. The default pattern is ".fcs", making sure that only fcs-files are selected. When you are already working with your data in *R*, it might be easier to use a *flowFrame* or *flowSet* from the *flowCore* package as input. This is also supported. If multiple paths or a *flowSet* are provided, all data will be concatenated. You should check and apply normalization if needed using other packages.

When reading the data, several preprocessing options are available. The data can be automatically compensated using a specified matrix, or using the `$SPILL` variable from the fcs-file (when `compensate` is `TRUE` but no value is given for `spillover`). The data can be transformed for specified columns. If no columns are provided, all columns from the spillover matrix will

be transformed. Finally, the data can be scaled. By default, it will scale to a mean of zero and standard deviation of one. However, specific scaling parameters can be set (see the base R `scale` function for more detail).

```
> set.seed(42)
> library(flowCore)
> library(FlowSOM)
> fileName <- system.file("extdata", "lymphocytes.fcs",
+                           package="FlowSOM")
> fSOM <- ReadInput(fileName, compensate = TRUE, transform = TRUE,
+                   toTransform=c(8:18), scale = TRUE)
> ff <- suppressWarnings(flowCore::read.FCS(fileName))
> fSOM <- ReadInput(ff, compensate = TRUE, transform = TRUE, scale = TRUE)
```

This function returns a FlowSOM object, which is actually a *list* containing several parameters. The data is stored as a matrix in `$data`, and all parameter settings to read the data are also stored. The begin and end indices of the subsets from the different files can be found in `$metadata`.

```
> str(fSOM, max.level = 2)

List of 12
 $ pattern      : chr ".fcs"
 $ compensate   : logi TRUE
 $ spillover    : num [1:11, 1:11] 1.00 4.84e-04 8.11e-04 8.63e-05 4.35e-04 ...
 .. attr(*, "dimnames")=List of 2
 $ transform    : logi TRUE
 $ toTransform  : chr [1:11] "FITC-A" "Pacific Blue-A" "AmCyan-A" "Qdot 605-A" ...
 $ transformFunction:Formal class 'transform' [package "flowCore"] with 2 slots
 $ scale        : logi TRUE
 $ prettyColnames : Named chr [1:18] "Time <Time>" "FSC-A <FSC-A>" "FSC-H <FSC-H>" "FSC-W <FSC-W>" ...
 .. attr(*, "names")= chr [1:18] "Time" "FSC-A" "FSC-H" "FSC-W" ...
 $ data         : num [1:19225, 1:18] -1.65 -1.65 -1.65 -1.65 -1.65 ...
 .. attr(*, "dimnames")=List of 2
 $ metaData     :List of 1
 ..$ /tmp/Rtmpo2NtZL/Rinst10d85bf2798b/FlowSOM/extdata/lymphocytes.fcs: num [1:2] 1 19225
 $ scaled.center : Named num [1:18] 3356 88594 68698 84405 36886 ...
 .. attr(*, "names")= chr [1:18] "Time" "FSC-A" "FSC-H" "FSC-W" ...
 $ scaled.scale  : Named num [1:18] 2038 15064 3236 12997 13923 ...
 .. attr(*, "names")= chr [1:18] "Time" "FSC-A" "FSC-H" "FSC-W" ...
 - attr(*, "class")= chr "FlowSOM"
```

### 3 Building the self-organizing map

The next step in the algorithm is to build a self-organizing map. Several parameters for the self-organizing map algorithm can be provided, such as the dimensions of the grid, the learning rate, the number of times the dataset has to be presented. However, the most important

parameter to decide is on which columns the self-organizing map should be trained. This should contain all the parameters that are useful to identify cell types, and exclude parameters of which you want to study the behavior on all cell types such as activation markers.

The `BuildSOM` function expects a FlowSOM object as input, and will return a FlowSOM object with all information about the self organizing map added in the `$map` parameter of the FlowSOM object.

```
> fSOM <- BuildSOM(fSOM,colsToUse = c(9,12,14:18))
> str(fSOM$map,max.level = 2)

List of 15
 $ xdim      : num 10
 $ ydim      : num 10
 $ rlen      : num 10
 $ mst       : num 1
 $ alpha     :List of 1
  ..$ : num [1:2] 0.05 0.01
 $ radius    :List of 1
  ..$ : num [1:2] 6 0
 $ init      : logi FALSE
 $ distf     : num 2
 $ grid      :'data.frame':      100 obs. of  2 variables:
  ..$ Var1: int [1:100] 1 2 3 4 5 6 7 8 9 10 ...
  ..$ Var2: int [1:100] 1 1 1 1 1 1 1 1 1 1 ...
  ..- attr(*, "out.attrs")=List of 2
 $ codes     : num [1:100, 1:7] -0.226 -0.397 -0.164 -0.352 -0.903 ...
  ..- attr(*, "dimnames")=List of 2
 $ mapping   : num [1:19225, 1:2] 10 99 2 96 42 92 91 43 81 27 ...
 $ nNodes    : int 100
 $ colsUsed  : num [1:7] 9 12 14 15 16 17 18
 $ medianValues: num [1:100, 1:18] 0.138 -0.643 0.256 0.18 0.302 ...
  ..- attr(*, "dimnames")=List of 2
 $ sdValues  : num [1:100, 1:18] 0.966 0.839 0.965 0.971 0.989 ...
  ..- attr(*, "dimnames")=List of 2
```

## 4 Building the minimal spanning tree

---

The third step of FlowSOM is to build the minimal spanning tree. This will again return a FlowSOM object, with extra information contained in the `$MST` parameter.

```
> fSOM <- BuildMST(fSOM,tSNE=TRUE)
> str(fSOM$MST)

List of 4
 $ graph:List of 10
  ..$ :List of 1
  .. ..$ 1: 'igraph.vs' Named int [1:2] 3 12
  .. ..- attr(*, "names")= chr [1:2] "3" "12"
  .. ..- attr(*, "env")=<weakref>
```

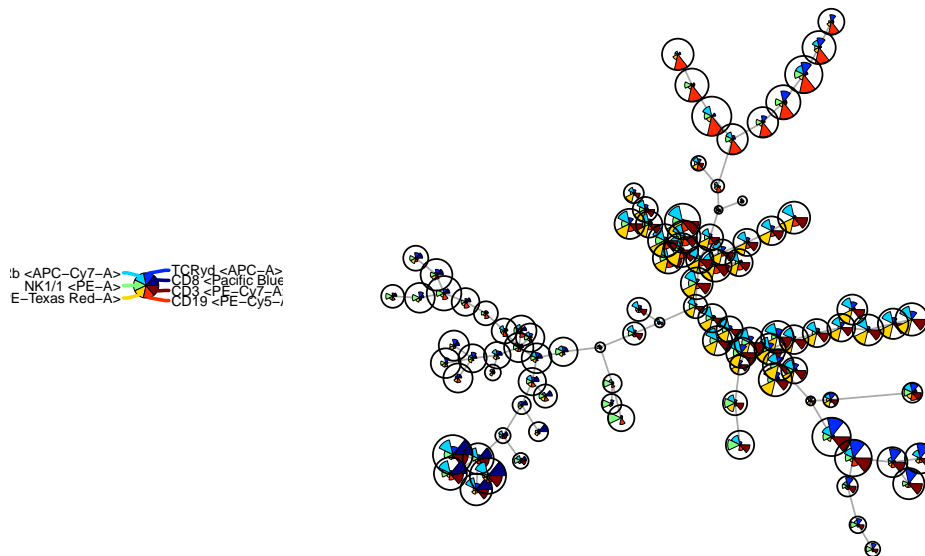
```

.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 2: 'igraph.vs' Named int [1:2] 3 11
.. .. ..- attr(*, "names")= chr [1:2] "3" "11"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 3: 'igraph.vs' Named int [1:2] 1 2
.. .. ..- attr(*, "names")= chr [1:2] "1" "2"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 4: 'igraph.vs' Named int 14
.. .. ..- attr(*, "names")= chr "14"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 5: 'igraph.vs' Named int 6
.. .. ..- attr(*, "names")= chr "6"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 6: 'igraph.vs' Named int [1:2] 5 16
.. .. ..- attr(*, "names")= chr [1:2] "5" "16"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 7: 'igraph.vs' Named int 17
.. .. ..- attr(*, "names")= chr "17"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 8: 'igraph.vs' Named int 9
.. .. ..- attr(*, "names")= chr "9"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 9: 'igraph.vs' Named int [1:2] 8 18
.. .. ..- attr(*, "names")= chr [1:2] "8" "18"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..$ :List of 1
.. ..$ 10: 'igraph.vs' Named int [1:2] 19 20
.. .. ..- attr(*, "names")= chr [1:2] "19" "20"
.. .. ..- attr(*, "env")=<weakref>
.. .. ..- attr(*, "graph")= chr "a5e9881d-d73d-4461-892b-5ca7e78b2f5f"
..- attr(*, "class")= chr "igraph"
$ l : num [1:100, 1:2] 1.68495 2.27903 2.02504 0.27503 0.00986 ...
$ l2 : num [1:100, 1:2] -129.5 -175.6 -177.3 106.6 41.5 ...
$ size : num [1:100] 13 12.62 14.19 5.63 13.52 ...

```

Once this step is finished, the FlowSOM object can be used for visualization. You can plot the nodes in several layouts ("MST": Minimal spanning tree (default), "grid": SOM grid, "tSNE": alternative layout, only possible when tSNE was TRUE in `BuildMST`)

```
> PlotStars(fSOM)
```

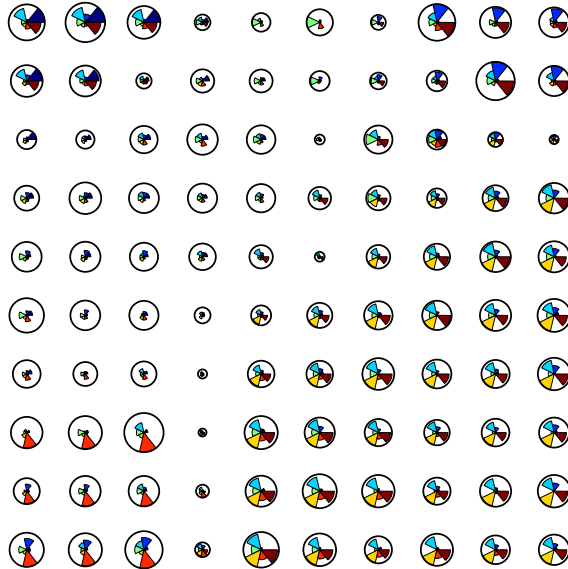


```
> PlotStars(fSOM, view="grid")
```

```

:b <APC-Cy7-A>
NK1/1 <PE-A>
E-Texas Red-A>
TCRvd1 <APC-A>
CD8 <Pacific Blue>
CD3 <PE-Cy7-A>
CD19 <PE-Cy5->

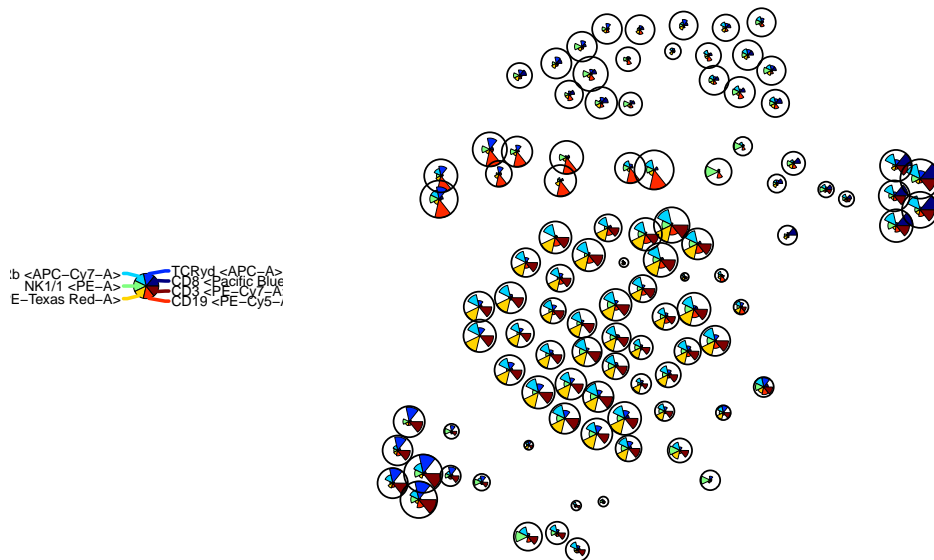
```



```

> PlotStars(fSOM, view="tSNE")

```

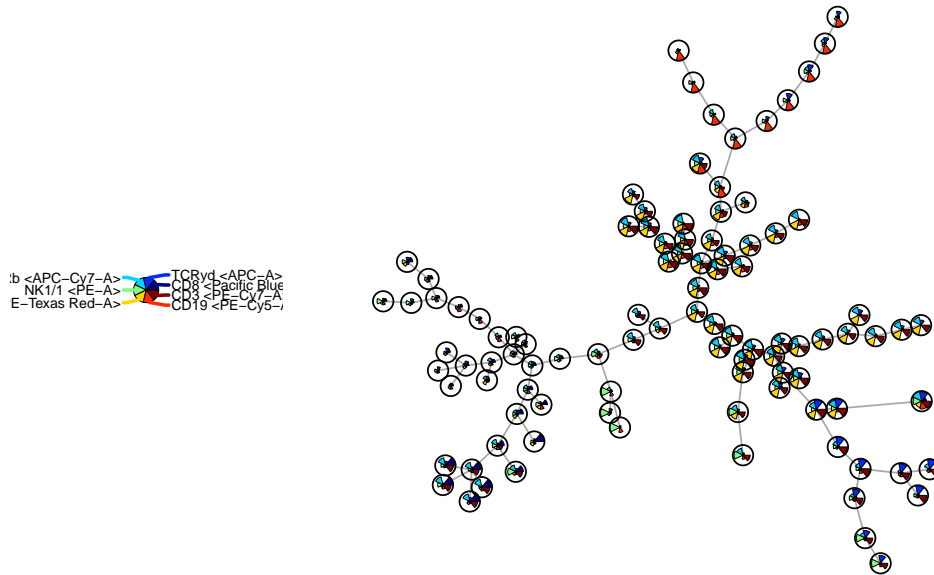


If you do not want the size to depend on the number of cells assigned to a node, you can reset the node size. This can be used in combination with any of the plotting functions.

```

> fSOM <- UpdateNodeSize(fSOM, reset=TRUE)
> fSOM$MST$size <- fSOM$MST$size/2
> PlotStars(fSOM)
> fSOM <- UpdateNodeSize(fSOM)
  
```

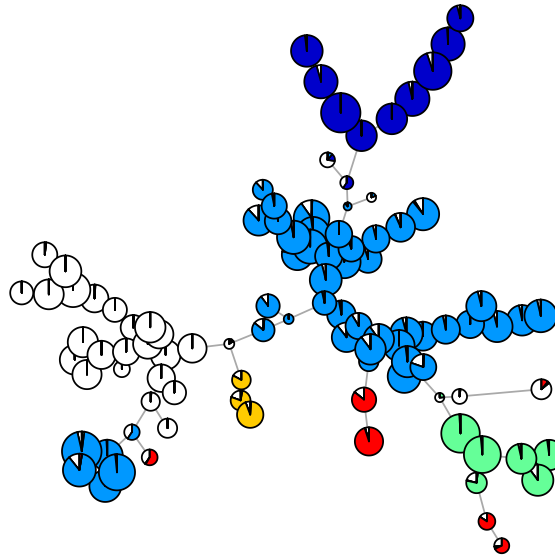




It might also be interesting to compare with a manual gating. The `cellTypes` can be any factor which has a value for each individual cell, so you can also map other clustering results.

```
> #<<>>=
> library(flowUtils)
> flowEnv <- new.env()
> ff_c <- compensate(ff,description(ff)$SPILL)
> colnames(ff_c)[8:18] <- paste("Comp-",colnames(ff_c)[8:18],sep="")
> gatingFile <- system.file("extdata","manualGating.xml",
+                             package="FlowSOM")
> gateIDs <- c( "B cells"=8,
+               "ab T cells"=10,
+               "yd T cells"=15,
+               "NK cells"=5,
+               "NKT cells"=6)
> cellTypes <- names(gateIDs)
> gatingResult <- ProcessGatingML(ff_c, gatingFile, gateIDs, cellTypes)
> PlotPies(fSOM,cellTypes=gatingResult$manual)
```

- Unknown
- B cells
- ab T cells
- yd T cells
- NK cells
- NKT cells



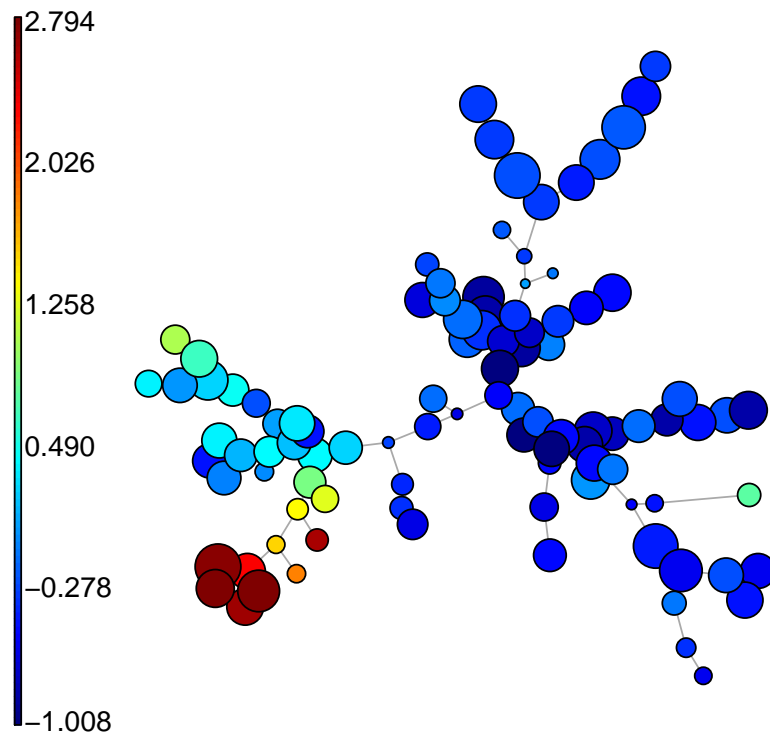
If you are interested in one specific marker, you can use the `PlotMarker` function.

```
> print(colnames(fsOM$map$medianValues))
```

\$P1N	\$P2N	\$P3N	\$P4N
"Time"	"FSC-A"	"FSC-H"	"FSC-W"
\$P5N	\$P6N	\$P7N	\$P8N
"SSC-A"	"SSC-H"	"SSC-W"	"FITC-A"
\$P9N	\$P10N	\$P11N	\$P12N
"Pacific Blue-A"	"AmCyan-A"	"Qdot 605-A"	"APC-A"
\$P13N	\$P14N	\$P15N	\$P16N
"Alexa Fluor 700-A"	"APC-Cy7-A"	"PE-A"	"PE-Texas Red-A"
\$P17N	\$P18N		
"PE-Cy5-A"	"PE-Cy7-A"		

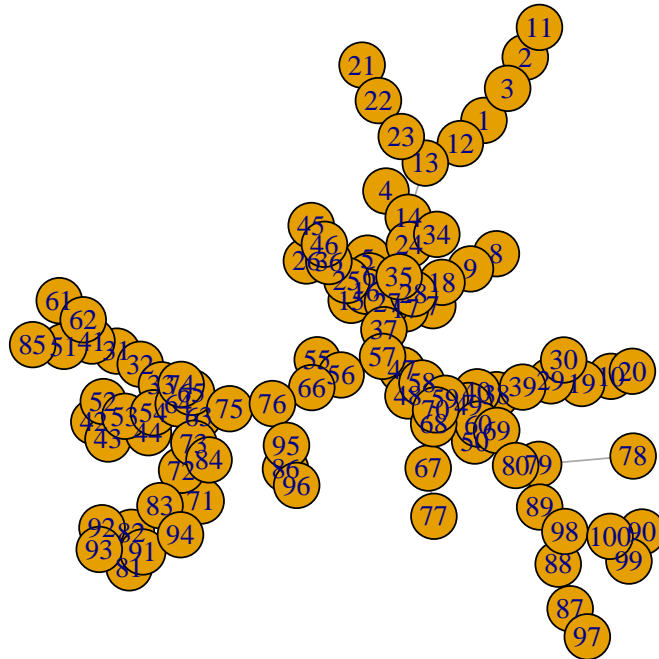
```
> PlotMarker(fsOM, "Pacific Blue-A")
```

## Pacific Blue-A



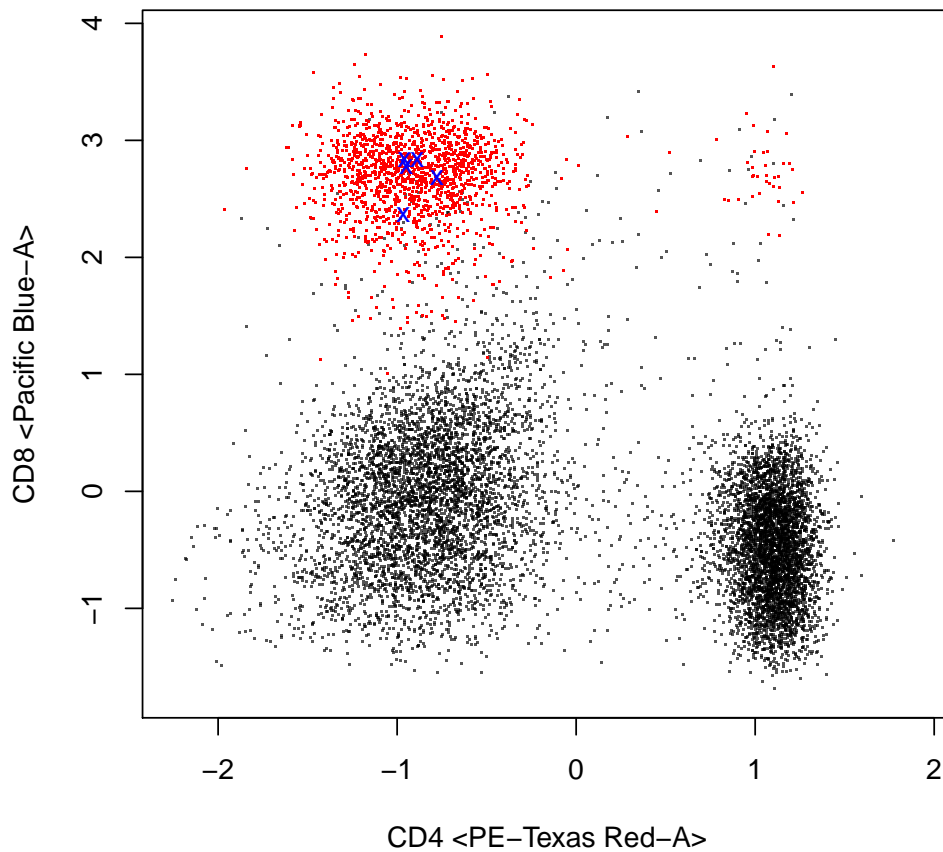
If you need to refer to the nodes, it might be useful to number them.

```
> PlotNumbers(UpdateNodeSize(fSOM, reset=TRUE))
```



You can use this number for a 2D scatter plot indicating the node values.

```
> PlotClusters2D(fSOM, "PE-Texas Red-A", "Pacific Blue-A", c(81, 82, 91, 92, 93))
```



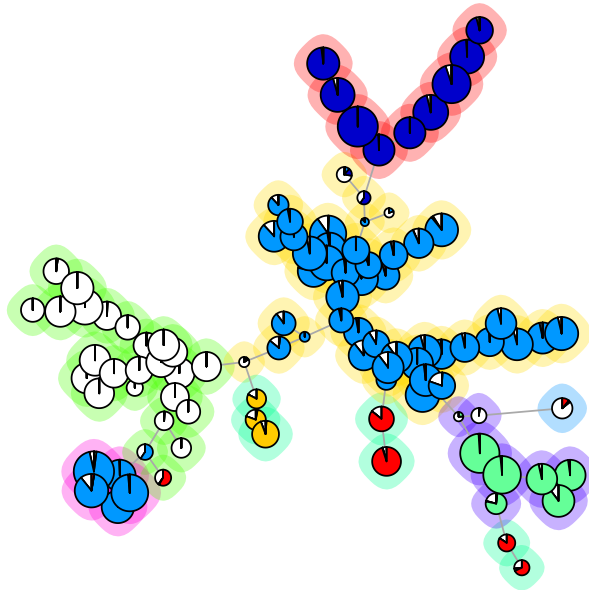
## 5 Meta-clustering the data

The fourth step of the FlowSOM algorithm is to perform a meta-clustering of the data. This can be the first step in further analysis of the data, and often gives a good approximation of manual gating results.

If you have background knowledge about the number of cell types you are looking for, it might be optimal to provide this number to the algorithm.

```
> #<<>>=
> metaClustering <- metaClustering_consensus(fSOM$map$codes,k=7)
> PlotPies(fSOM,cellTypes=gatingResult$manual,
+         backgroundValues = as.factor(metaClustering))
```

- Unknown
- B cells
- ab T cells
- yd T cells
- NK cells
- NKT cells



Background

- 1
- 2
- 3
- 4
- 5
- 6
- 7

You can also extract the meta-clustering for each cell individually

```
> metaClustering_perCell <- metaClustering[fSOM$map$mapping[,1]]
```

## 6 Detecting nodes with a specific pattern

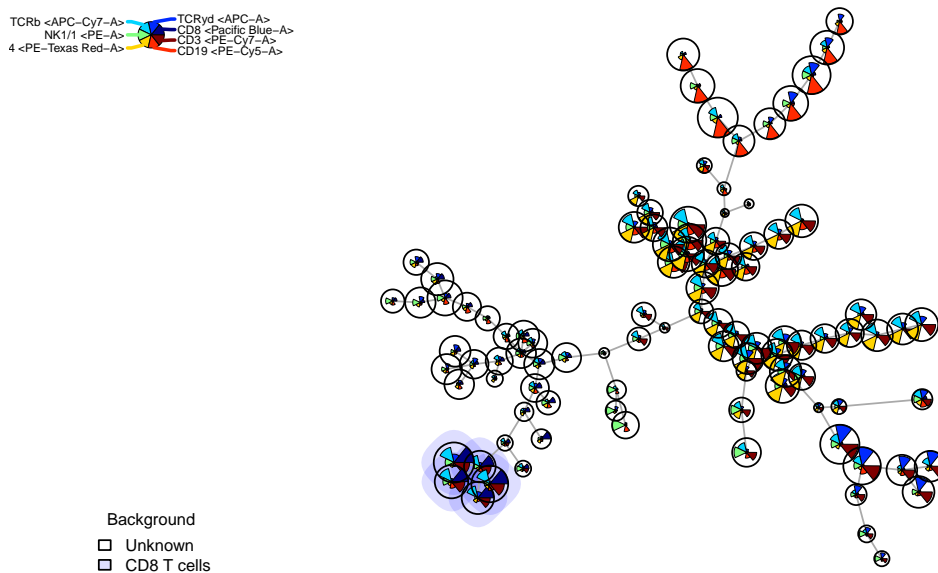
If you do not have a manual gating to map on the tree, it might be time-consuming to interpret all the different nodes. Therefore, you can also query the tree to indicate nodes similar to a specified pattern. This function is still being optimized, so make sure to check the marker values to see if it corresponds to your expectations.

```
> # Look for CD8+ ab T cells
> query <- c("PE-Cy7-A" = "high", #CD3
+          "APC-Cy7-A" = "high", #TCRb
+          "Pacific Blue-A" = "high") #CD8
> query_res <- QueryStarPlot(UpdateNodeSize(fSOM, reset=TRUE), query,
+                             plot = FALSE)
```

```

> cellTypes <- factor(rep("Unknown",49),levels=c("Unknown","CD8 T cells"))
> cellTypes[query_res$selected] <- "CD8 T cells"
> PlotStars(fSOM,
+           backgroundValues=cellTypes,
+           backgroundColor=c("#FFFFFF00", "#0000FF22"))

```



## 7 Comparing different groups

It is possible to compare between groups with the FlowSOM package as well. The tree should be built on either a concatenation of all files, or a representative subset of all cell types. Then a list identifying which files belong to specific groups should be defined, and the differences will be computed. For a smaller number of samples, you can look at the fold change between the groups. For coloring, a threshold is used. A threshold of 0.50 means the difference should be at least 50% of the max of both groups, which corresponds with a 2-fold change. The

higher the threshold, the stricter, a threshold of 0 will colour each node. For a larger number of samples you can also use a wilcox test. This will be selected when a value is provided for the `p_tresh` parameter.

```

> library(FlowSOM)
> # Build the FlowSOM tree on the example file
> fileName <- system.file("extdata", "lymphocytes.fcs", package="FlowSOM")
> flowSOM.res <- FlowSOM(fileName, compensate=TRUE, transform=TRUE,
+ scale=TRUE, colsToUse=c(9,12,14:18), nClus = 10, seed=1)
> # Have a look at the resulting tree
> # PlotStars(flowSOM.res[[1]], backgroundValues = as.factor(flowSOM.res[[2]]))
>
> # Select all cells except the branch that corresponds with automated
> # cluster 7 (CD3+ TCRyd +) and write to another file for the example
> # In practice you would not generate any new file but
> # use your different files from your different groups
> ff <- flowCore::read.FCS(fileName)
> ff_tmp <- ff[flowSOM.res[[1]]$map$mapping[,1] %in% which(flowSOM.res[[2]] != 7),]
> flowCore::write.FCS(ff_tmp, file="ff_tmp.fcs")

[1] "ff_tmp.fcs"

> # Make an additional file without cluster 7 and double amount of cluster 10
> ff_tmp <- ff[c(which(flowSOM.res[[1]]$map$mapping[,1] %in% which(flowSOM.res[[2]] != 7)),
+ which(flowSOM.res[[1]]$map$mapping[,1] %in% which(flowSOM.res[[2]] == 5))),]
> flowCore::write.FCS(ff_tmp, file="ff_tmp2.fcs")

[1] "ff_tmp2.fcs"

> # Compare the original file with the two new files we made
> groupRes <- CountGroups(flowSOM.res[[1]],
+ groups=list("AllCells"=c(fileName),
+ "Without_ydTcells"=c("ff_tmp.fcs", "ff_tmp2.fcs")))

[1] "/tmp/Rtmpo2NtZL/Rinst10d85bf2798b/FlowSOM/extdata/lymphocytes.fcs"
[1] "ff_tmp.fcs"
[1] "ff_tmp2.fcs"

> # PlotGroups(flowSOM.res[[1]], groupRes)
>
> # Compare only the file with the double amount of cluster 10
> groupRes <- CountGroups(flowSOM.res[[1]],
+ groups=list("AllCells"=c(fileName),
+ "Without_ydTcells"=c("ff_tmp2.fcs")))

[1] "/tmp/Rtmpo2NtZL/Rinst10d85bf2798b/FlowSOM/extdata/lymphocytes.fcs"
[1] "ff_tmp2.fcs"

> PlotGroups(flowSOM.res[[1]], groupRes)
$Without_ydTcells
[1] -- -- -- --
[5] -- -- -- --
[9] -- -- -- --
[13] -- -- -- --
[17] -- -- -- --

```

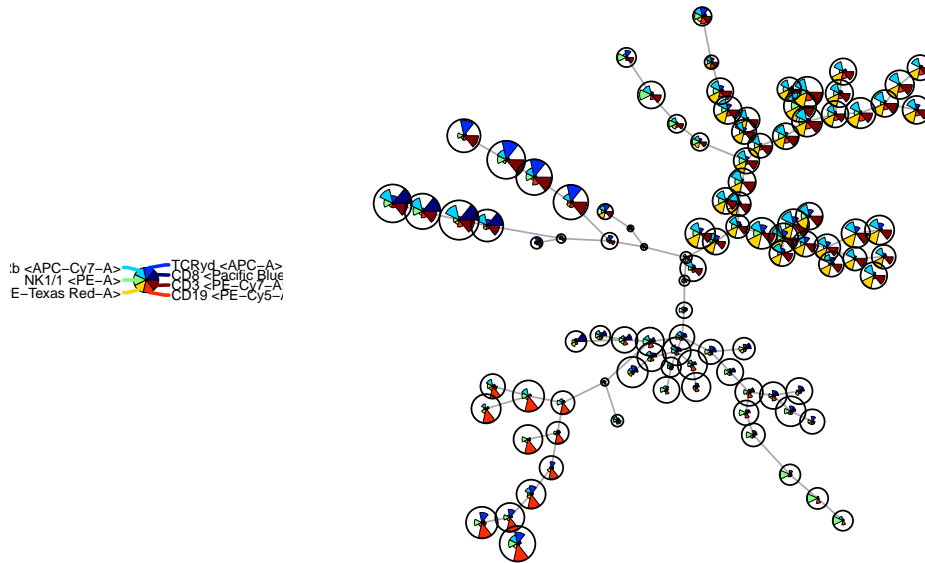


```

[21] --          --          --          --
[25] --          --          --          --
[29] --          --          --          --
[33] --          --          --          --
[37] --          --          --          --
[41] --          --          --          --
[45] --          --          --          Without_ydTcells
[49] --          --          AllCells    AllCells
[53] --          --          --          --
[57] Without_ydTcells Without_ydTcells --          --
[61] AllCells        AllCells        --          --
[65] --          --          --          --
[69] --          --          --          --
[73] --          --          --          --
[77] --          --          --          --
[81] --          --          --          --
[85] --          --          --          --
[89] --          --          --          --
[93] --          --          --          --
[97] --          --          --          --
Levels: -- AllCells Without_ydTcells

```

tZL/Rinst10d85bf2798b/FlowSOM/extdata/ly



## 8 Summary

In summary, the FlowSOM package provides some new ways to look at cytometry data. It can help to keep an overview of how all markers are behaving on different cell types, and to reduce the probability of overlooking interesting things that are present in the data.