

# Multiome tutorial - MultiAssayExperiment

**Xiaosai Yao**

Nov 3rd, 2022

## Package

epiregulon 1.0.16

## Contents

1	Introduction	2
2	Installation	2
3	Data preparation	2
4	Quick start	3
4.1	Retrieve bulk TF ChIP-seq binding sites	3
4.2	Link ATAC-seq peaks to target genes	4
4.3	Add TF motif binding to peaks	4
4.4	Generate regulons	5
4.5	Network pruning (highly recommended)	5
4.6	Add Weights	6
4.7	Calculate TF activity	7
4.8	Perform differential activity	7
4.9	Visualize the results	7
4.10	Geneset enrichment	10
4.11	Network analysis	11
4.12	Differential networks	11
5	Session Info	14

## 1 Introduction

This tutorial walks through the same dataset used in the “multiome tutorial - archR workflow”. This is a dataset generated by infecting LNCaP cells with NKX2-1 and GATA6 to examine the effects of these TFs on AR activity.

## 2 Installation

Epiregulon is currently available on R/dev

```
library(epiregulon)
```

Alternatively, you could install from gitlab

```
devtools::install_gitlab(repo='scwg/gene-transcriptional-network/activity-inference/Epiregulon',
                        auth_token = "<gitlab token>",
                        host = "https://code.roche.com/" )

library(epiregulon)
```

## 3 Data preparation

Single cell preprocessing needs to be performed by user’s favorite methods prior to using Epiregulon. The following components are required: 1. Peak matrix from scATAC-seq 2. Gene expression matrix from either paired or unpaired scRNA-seq. RNA-seq integration needs to be performed for unpaired dataset. 3. Dimensionality reduction matrix from with either single modalities or joint scRNA-seq and scATAC-seq

Multiome data can now be conveniently processed by `initiate.archr` and then `gp.sa.archr` to obtain peak matrices. Finally, the archR project can be uploaded into DatasetDB as a `MultiAssayExperiment` object using `maw.archr::importArchr` or `maw.archr::create.mae.with.multiple.scs.from.archr`

```
# load the MAE object
library(SingleCellExperiment)
mae <- dsassembly::getDataset("DS000013080")
#> 'version=' not specified, using the latest version (2) instead
#> Error in read_csv(path, is_compressed = identical(compression, "gzip"), :
#>   encountered empty line in a file with non-zero columns
#>   falling back to 'read.csv' for a malformed CSV data frame

# peak matrix
PeakMatrix <- mae[["PeakMatrix"]]
rownames(PeakMatrix) <- rowData(PeakMatrix)$idx

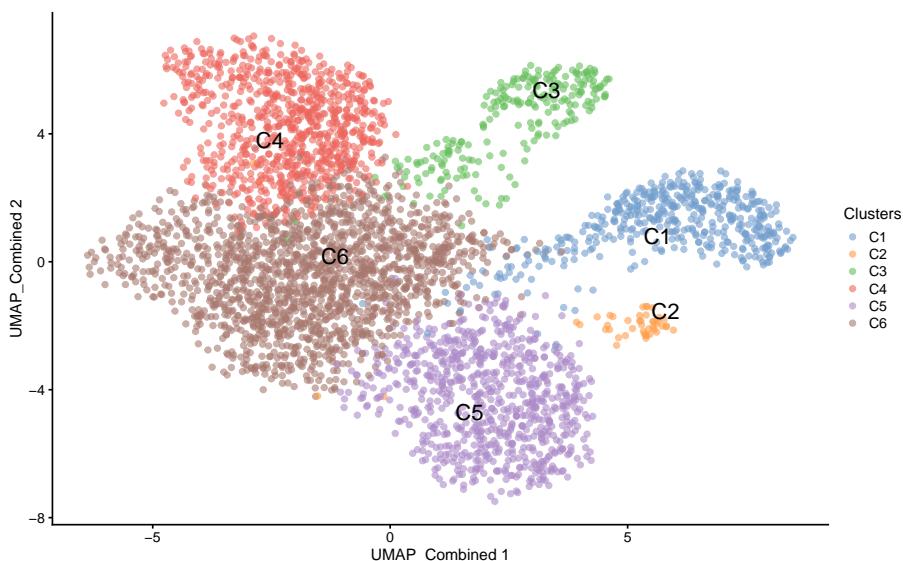
# expression matrix
GeneExpressionMatrix <- mae[["GeneExpressionMatrix"]]
rownames(GeneExpressionMatrix) <- rowData(GeneExpressionMatrix)$name
```

## Multiome tutorial - MultiAssayExperiment

```
# dimensional reduction matrix
reducedDimMatrix <- reducedDim(mae[['TileMatrix500']], "LSI_ATAC")
```

Visualize singleCellExperiment by UMAP

```
# transfer UMAP_combined from TileMatrix to GeneExpressionMatrix
reducedDim(GeneExpressionMatrix, "UMAP_Combined") <- reducedDim(mae[['TileMatrix500']], "UMAP_Combined")
scater::plotReducedDim(GeneExpressionMatrix,
                       dimred = "UMAP_Combined",
                       text_by = "Clusters",
                       colour_by = "Clusters")
```



## 4 Quick start

### 4.1 Retrieve bulk TF ChIP-seq binding sites

First, we retrieve the information of TF binding sites collected from Cistrome and ENCODE ChIP-seq, which are hosted on Genomitory. Currently, human genomes HG19 and HG38 and mouse mm10 are available.

```
grl <- getTFMotifInfo(genome = "hg38")
#> redirecting from 'GMTY162:hg38_motif_bed_granges@REVISION-4' to 'GMTY162:hg38_motif_bed_granges@24c22e4f42'
head(grl)
#> GRangesList object of length 6:
#> $`5-hmC`
#> GRanges object with 24048 ranges and 0 metadata columns:
#>   seqnames      ranges strand
#>   <Rle>      <IRanges>  <Rle>
#>   [1]    chr1    10000-10685     *
#>   [2]    chr1    13362-13694     *
```

## Multiome tutorial - MultiAssayExperiment

```
#>      [3] chr1    29631-29989      *
#>      [4] chr1    40454-40754      *
#>      [5] chr1    135395-135871     *
#>      ...
#>      ...      ...      ...
#> [24044] chrY    56864377-56864627     *
#> [24045] chrY    56876124-56876182     *
#> [24046] chrM    84-2450      *
#> [24047] chrM    13613-14955     *
#> [24048] chrM    15134-16490      *
#> -----
#> seqinfo: 25 sequences from an unspecified genome; no seqlengths
#>
#> ...
#> <5 more elements>
```

## 4.2 Link ATAC-seq peaks to target genes

Next, we compute peak to gene correlations using a custom algorithm that has similar performance to ArchR's P2G function.

```
p2g <- calculateP2G(peakMatrix = PeakMatrix,
                      expMatrix = GeneExpressionMatrix,
                      reducedDim = reducedDimMatrix)
#> Using epiregulon to compute peak to gene links...
#> performing k means clustering to form metacells
#> Computing correlation

head(p2g)
#>   idxATAC chr start   end idxRNA      target Correlation distance      FDR
#> 67       1 chr1 817121 817621      19   FAM41C  0.6869051  50578 0.07194023
#> 478      9 chr1 920452 920952      29 AL645608.7  0.5527224  75097 0.33240065
#> 76       10 chr1 920987 921487     19   FAM41C  0.5870781  50587 0.24966893
#> 88       22 chr1 960317 960817     19   FAM41C  0.6859528  89917 0.07303724
#> 445      22 chr1 960317 960817     28   PERM1   0.6897232  18210 0.06877207
#> 691      22 chr1 960317 960817     33   AGRN   0.7094602  57301 0.04868069
```

## 4.3 Add TF motif binding to peaks

The next step is to add the TF binding information by overlapping regions of the peak matrix with the bulk chip-seq database loaded in 2. The user can supply either an archR project path and this function will retrieve the peak matrix, or a peakMatrix in the form of a Granges object or RangedSummarizedExperiment.

```
overlap <- addTFMotifInfo(grl = grl, p2g = p2g, peakMatrix = PeakMatrix)
#> Computing overlap...
#> Success!
head(overlap)
```

```
#>   idxATAC idxTF      tf
#> 1       1     2 5-mC
#> 2       1    22 AML1-ETO
#> 3       1    25 AR
#> 4       1    49 ATF1
#> 5       1    50 ATF2
#> 6       1    51 ATF3
```

## 4.4 Generate regulons

A long format dataframe, representing the inferred regulons, is then generated. The dataframe consists of three columns:

- tf (transcription factor)
- target gene
- peak to gene correlation between tf and target gene

```
regulon <- getRegulon(p2g = p2g, overlap = overlap, aggregate = FALSE)
head(regulon)
#>   idxATAC idxTF      tf chr start   end idxRNA target      corr distance      FDR
#> 1       1     2 5-mC chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
#> 2       1    22 AML1-ETO chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
#> 3       1    25      AR chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
#> 4       1    49      ATF1 chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
#> 5       1    50      ATF2 chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
#> 6       1    51      ATF3 chr1 817121 817621    19 FAM41C 0.6869051 50578 0.07194023
```

## 4.5 Network pruning (highly recommended)

Epiregulon prunes the network by performing tests of independence on the observed number of cells jointly expressing transcription factor (TF), regulatory element (RE) and target gene (TG) vs the expected number of cells if TF/RE and TG are independently expressed. We implement two tests, the binomial test and the chi-square test. In the binomial test, the expected probability is  $P(\text{TF}, \text{RE}) * P(\text{TG})$ , and the number of trials is the total number of cells, and the observed successes is the number of cells jointly expressing all three elements. In the chi-square test, the expected probability for having all 3 elements active is also  $P(\text{TF}, \text{RE}) * P(\text{TG})$  and the probability otherwise is  $1 - P(\text{TF}, \text{RE}) * P(\text{TG})$ . The observed cell count for the active category is the number of cells jointly expressing all three elements, and the cell count for the inactive category is  $n - n_{\text{triple}}$ .

We calculate cluster-specific p-values if users supply cluster labels. This is useful if we are interested in cluster-specific networks. The pruned regulons can then be used to visualize differential networks for transcription factors of interest. See section on differential networks.

```
pruned.regulon <- pruneRegulon(expMatrix = GeneExpressionMatrix,
                                exp_assay = "counts",
                                peakMatrix = PeakMatrix,
```

```
peak_assay = "counts",
test = "binom",
regulon = regulon[regulon$tf %in% c("NKX2-1", "GATA6", "FOXA1", "FOXA2", "AR"),
#regulon = regulon,
clusters = GeneExpressionMatrix$Clusters,
prune_value = "pval",
regulon_cutoff = 0.05,
BPPARAM = BiocParallel::MulticoreParam(workers = 2, progressbar = TRUE)
)
#> pruning network with binom tests using a regulon cutoff of pval<0.05
```

## 4.6 Add Weights

While the 'pruneRegulon' function provides statistics on the joint occurrence of TF-RE-TG, we would like to further estimate the strength of regulation. Biologically, this can be interpreted as the magnitude of gene expression changes induced by transcription factor activity. Epiregulon estimates the regulatory potential using one of the four measures: 1) correlation between TF and target gene expression, 2) mutual information between the TF and target gene expression, 3) Wilcoxon test statistics of target gene expression in cells jointly expressing all 3 elements vs cells that do not, or 4) log 2 fold difference of target gene expression in cells jointly expressing all 3 elements vs cells that do not.

Three measures (correlation, Wilcoxon statistics and log 2 fold difference) give both the magnitude and directionality of changes whereas mutational information is always positive. The correlation and mutual information statistics are computed on the grouped pseudobulks by user-supplied cluster labels, whereas the Wilcoxon and log fold change group cells based on the joint expression of TF, RE and TG in each single cell.

```
regulon.w <- addWeights(regulon = pruned.regulon,
                         sce = GeneExpressionMatrix,
                         exprs_values = "counts",
                         peakMatrix = PeakMatrix,
                         peak_assay = "counts",
                         cluster_factor = "Clusters",
                         block_factor = NULL,
                         method = "corr")
```

```

head( regulon.w )
#>      idxATAC idxTF tf   chr     start       end idxRNA target      corr distance          FDR      pval
#> 1256430    46808    25 AR chr16 70289316 70289816 13581 AARS 0.5743454        0 0.279160621 0.038119
#> 639674     24880    25 AR chr11 106121522 106122022 6438 AASDHPPPT 0.6784558 43661 0.082155624 0.187058
#> 1504723    53005    25 AR chr17 81038117 81038617 15850 AATK 0.5225027 91329 0.407852377 0.015859
#> 2004650    69758    25 AR chr20 31698801 31699301 21689 ABALON 0.8249891 20204 0.002393355 0.002463
#> 2005460    69759    25 AR chr20 31706227 31706727 21689 ABALON 0.7841009 12778 0.008905490 0.015938
#> 1123276    44160    25 AR chr16 2339633 2340133 12490 ABCA3 0.6465149 9710 0.129436606 0.211675
#>      stats_C5 stats_C6 padj_all padj_C1 padj_C2 padj_C3 padj_C4 padj_C5 padj_C6      weight
#> 1256430 1.8704972 0.1816975        1       1       1       1       1       1 0.7300760
#> 639674 0.0000000 2.0912769        1       1       1       1       1       1 -0.4385840
#> 1504723 1.4713907 2.0724933        1       1       1       1       1       1 0.6717650

```

## Multiome tutorial - MultiAssayExperiment

```
#> 2004650 1.1915363 0.6829362      1      1      1      1      1      1      1      1 0.1362656
#> 2005460 0.0000000 1.1640070      1      1      1      1      1      1      1      1 0.1362656
#> 1123276 -0.6283722 -0.2432893     1      1      1      1      1      1      1      1 0.6317828
```

## 4.7 Calculate TF activity

Finally, the activities for a specific TF in each cell are computed by averaging expressions of target genes linked to the TF weighted by the test statistics of choice, chosen from either correlation, mutual information, Wilcoxon test statistics or log fold change.

$$y = \frac{1}{n} \sum_{i=1}^n x_i * weights_i$$

where  $y$  is the activity of a TF for a cell  $n$  is the total number of targets for a TF  $x_i$  is the log count expression of target i where i in  $\{1,2,\dots,n\}$   $weights_i$  is the weight of TF and target i

```
score.combine <- calculateActivity(sce = GeneExpressionMatrix,
                                      regulon = regulon.w,
                                      mode = "weight",
                                      method = "weightedMean",
                                      assay = "counts")
```

## 4.8 Perform differential activity

```
markers <- findDifferentialActivity(activity_matrix = score.combine,
                                         groups = GeneExpressionMatrix$Clusters,
                                         pval.type = "some",
                                         direction = "up",
                                         test.type = "t")
```

Take the top TFs

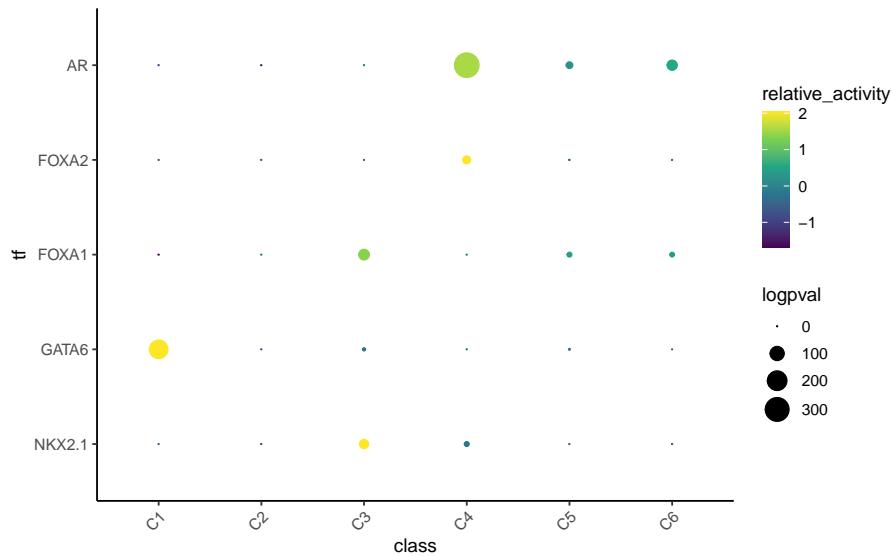
```
markers.sig <- getSigGenes(markers, topgenes = 5 )
#> Using a logFC cutoff of 0.8 for class C1
#> Using a logFC cutoff of 0 for class C2
#> Using a logFC cutoff of 0.3 for class C3
#> Using a logFC cutoff of 0.1 for class C4
#> Using a logFC cutoff of 0.1 for class C5
#> Using a logFC cutoff of 0.1 for class C6
```

## 4.9 Visualize the results

First visualize the known differential TFs by bubble plot

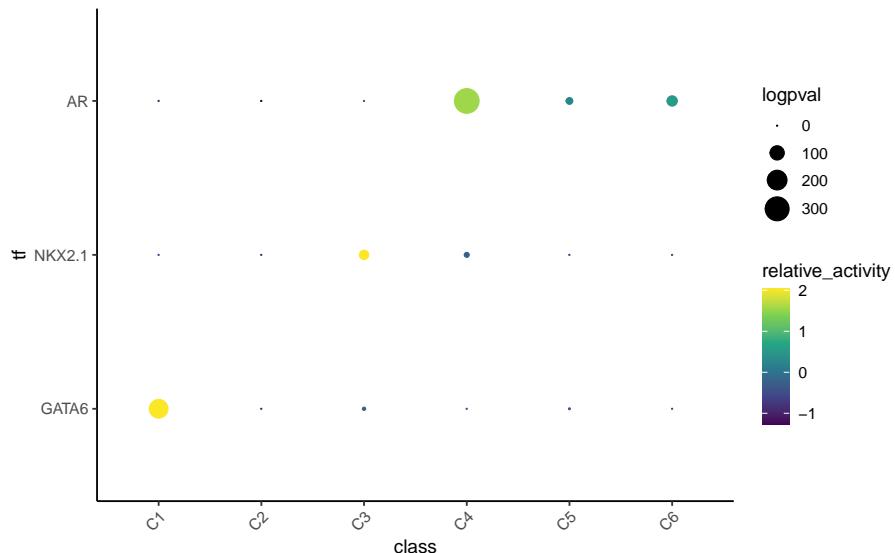
## Multiome tutorial - MultiAssayExperiment

```
plotBubble(activity_matrix = score.combine,
           tf = c("NKX2-1", "GATA6", "FOXA1", "FOXA2", "AR"),
           class = GeneExpressionMatrix$Clusters)
```



Then visualize the most differential TFs by clusters

```
plotBubble(activity_matrix = score.combine,
           tf = markers.sig$tf,
           class = GeneExpressionMatrix$Clusters)
```

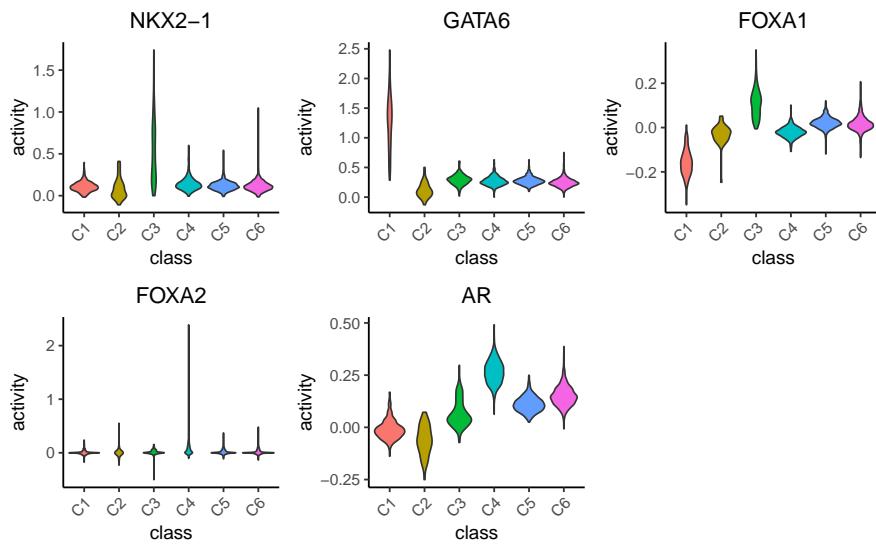


Visualize the known differential TFs by violin plot. Note there is no activity calculated for SOX2 because the expression of SOX2 is 0 in all cells.

```
plotActivityViolin(activity_matrix = score.combine,
                   tf = c("NKX2-1", "GATA6", "FOXA1", "FOXA2", "AR", "SOX2"))
```

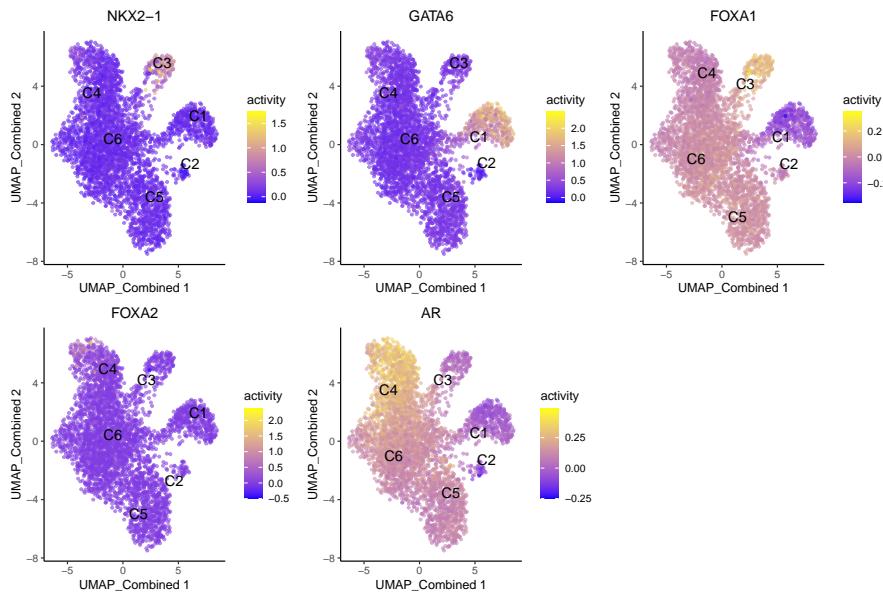
## Multiome tutorial - MultiAssayExperiment

```
class = GeneExpressionMatrix$Clusters)
#> SOX2 not found in activity matrix. Excluded from plots
```



Visualize the known differential TFs by UMAP

```
plotActivityDim(sce = GeneExpressionMatrix, activity_matrix=score.combine,
                 tf = c("NKX2-1", "GATA6", "FOXA1", "FOXA2", "AR", "SOX2"),
                 dimtype = "UMAP_Combined",
                 label = "Clusters",
                 point_size=1,
                 ncol = 3)
#> SOX2 not found in activity matrix. Excluded from plots
```



## 4.10 Geneset enrichment

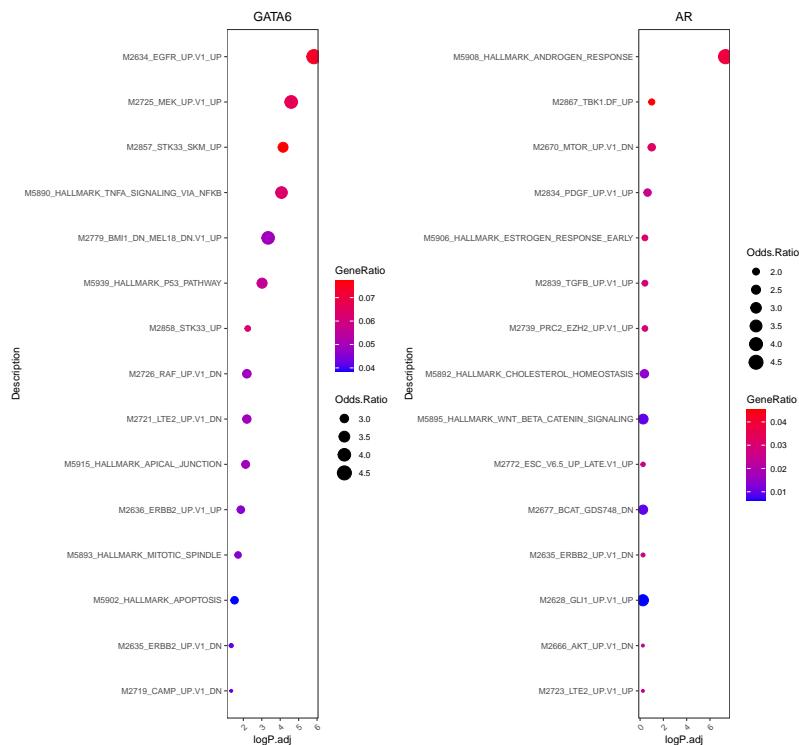
Sometimes we are interested to know what pathways are enriched in the regulon of a particular TF. We can perform geneset enrichment using the enricher function from [clusterProfiler](#).

```
#retrieve genesets
H <- EnrichmentBrowser::getGenesets(org = "hsa", db = "msigdb", cat = "H",
                                         gene.id.type = "SYMBOL")
#> Using cached version from 2022-11-02 02:21:28
C6 <- EnrichmentBrowser::getGenesets(org = "hsa", db = "msigdb", cat = "C6",
                                         gene.id.type = "SYMBOL")
#> Using cached version from 2022-11-02 02:21:34

#combine genesets and convert genesets to be compatible with enricher
gs <- c(H,C6)
gs.list <- do.call(rbind,lapply(names(gs), function(x)
  {data.frame(gs=x, genes=gs[[x]])}))

enrichresults <- regulonEnrich(TF = c("GATA6","AR"),
                                 regulon = regulon.w,
                                 corr = "weight",
                                 corr_cutoff = 0.1,
                                 genesets = gs.list)

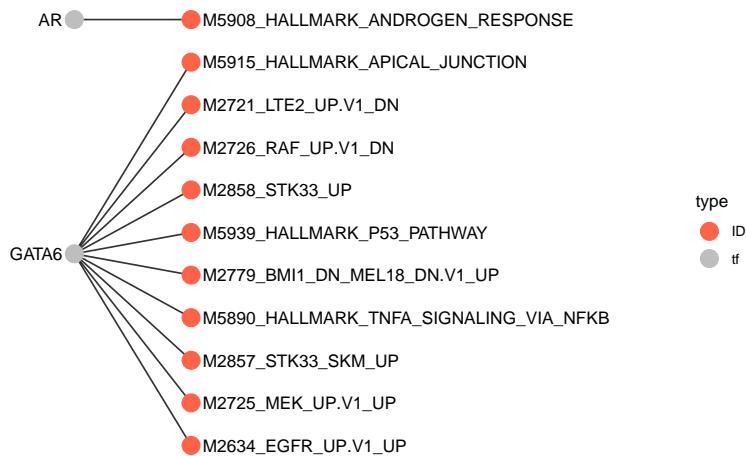
#plot results
enrichPlot(results = enrichresults, )
```



## 4.11 Network analysis

We can visualize the genesets as a network

```
plotGseaNetwork(tf = names(enrichresults),
                enrichresults = enrichresults,
                p.adj_cutoff = 0.1,
                ntop_pathways = 10)
```

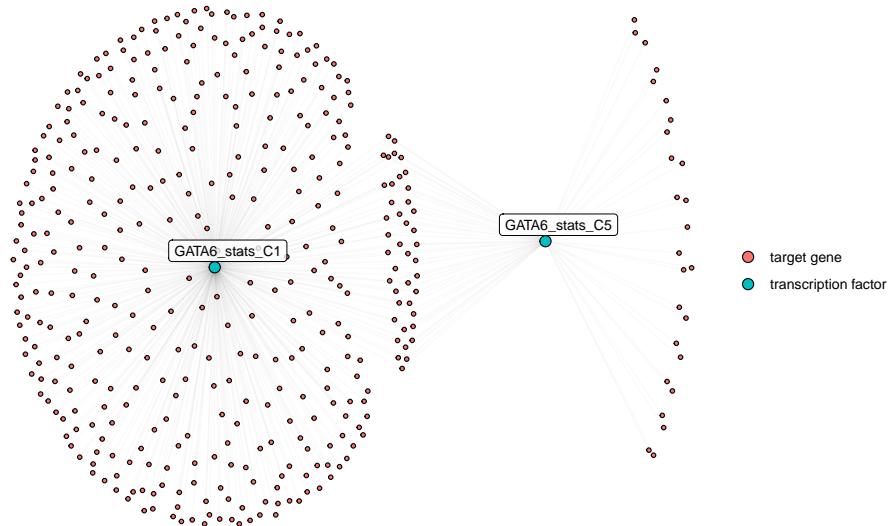


## 4.12 Differential networks

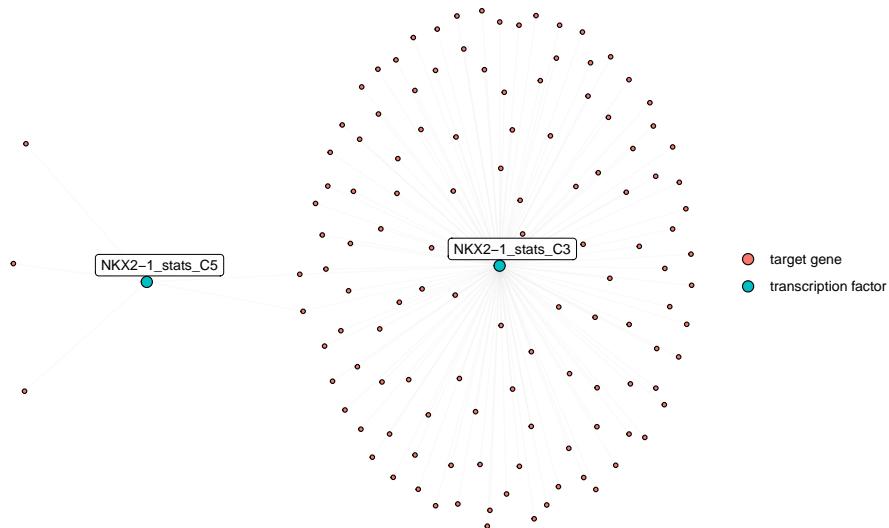
We are interested in understanding the differential networks between two conditions and determining which transcription factors account for the differences in the topology of networks. The pruned regulons with cluster-specific test statistics computed by `pruneRegulon` can be used to generate cluster-specific networks based on user-defined cutoffs and to visualize differential networks for transcription factors of interest. In this dataset, the GATA6 gene was only expressed in cluster 1 (C1) and NKX2-1 was only expressed in cluster 3 (C3). If we visualize the target genes of GATA6, we can see that C1 has many more target genes of GATA6 compared to C5, a cluster that does not express GATA6. Similarly, NKX2-1 target genes are confined to C3 which is the only cluster that exogenously expresses NKX2-1.

```
plotDiffNetwork(pruned.regulon,
                cutoff = 1,
                tf = c("GATA6"),
                groups = c("stats_C1", "stats_C5"),
                layout = "stress")
```

## Multiome tutorial - MultiAssayExperiment



```
plotDiffNetwork(pruned.regulon,
                cutoff = 1,
                tf = c("NKX2-1"),
                groups = c("stats_C3", "stats_C5"),
                layout = "stress")
```



We can also visualize how transcription factors relate to other transcription factors in each cluster.

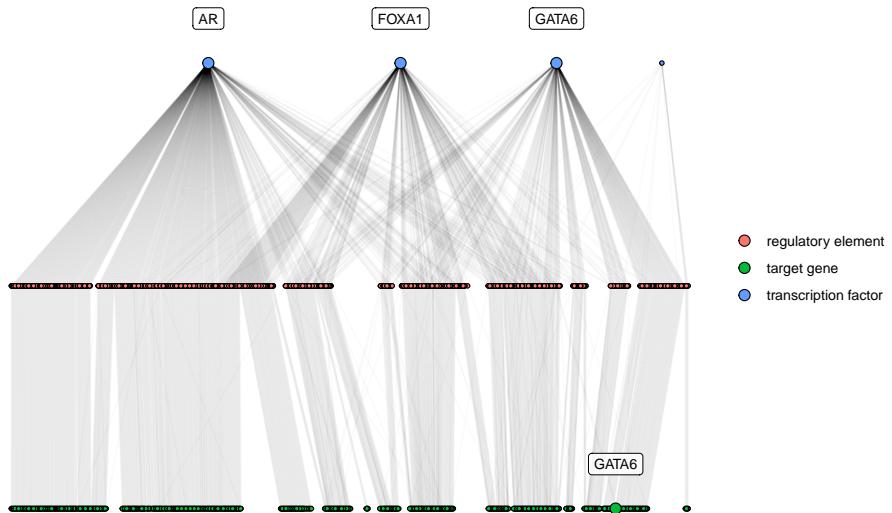
```
C1_network <- buildGraph(pruned.regulon[pruned.regulon$stats_C1>1,], weights = "stats_C1")
C5_network <- buildGraph(pruned.regulon[pruned.regulon$stats_C5>1,], weights = "stats_C5")

plotEpiRegulonNetwork(C1_network,
                      layout = "sugiyama",
```

## Multiome tutorial - MultiAssayExperiment

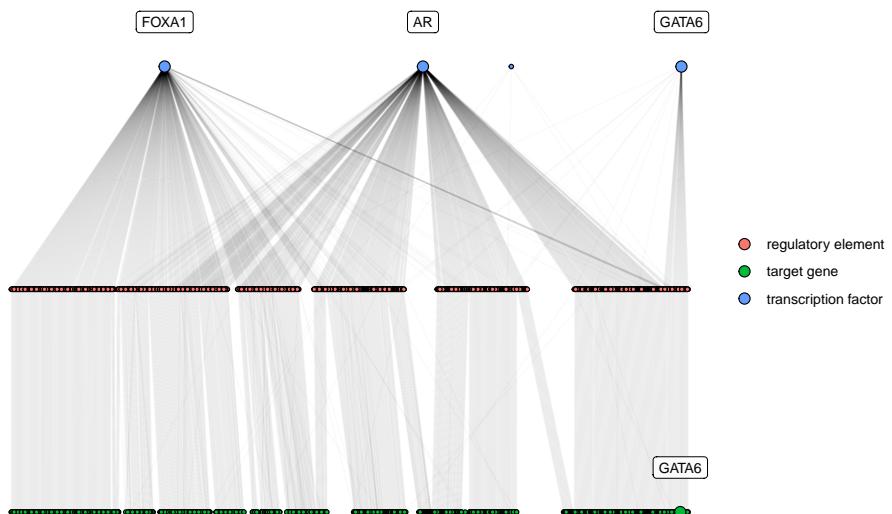
```
tf_to_highlight = c("GATA6", "FOXA1", "AR") + ggplot2::ggtitle ("C1")
```

C1



```
plotEpiregulonNetwork(C5_network,
                      layout = "sugiyama",
                      tf_to_highlight = c("GATA6", "FOXA1", "AR")) + ggplot2::ggtitle ("C5")
```

C5



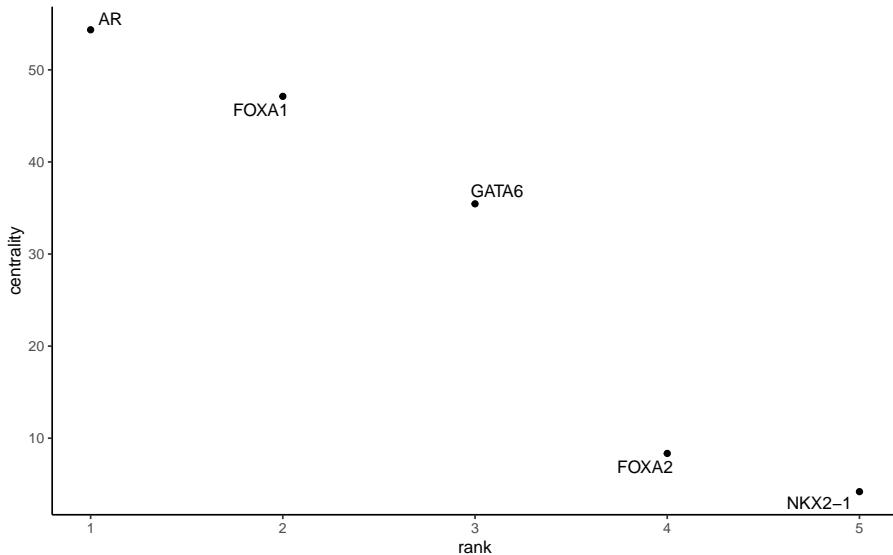
To systematically examine the differential network topology between two clusters, we perform an edge subtraction between two graphs, using weights computed by `pruneRegulon`. We then calculate the degree centrality of the weighted differential graphs and if desired, normalize the differential centrality against the total number of edges. The default normalization function is `sqrt` as it preserves both the difference in the number of edges (but scaled by `sqrt`) and the differences in the weights. If the user only wants to examine the differences in the averaged weights, the `FUN` argument can be changed to `identity`. Finally, we rank the transcription factors by (normalized) differential centrality.

## Multiome tutorial - MultiAssayExperiment

```
# rank by differential centrality
C1_network <- buildGraph(pruned.regulon, weights = "stats_C1")
C5_network <- buildGraph(pruned.regulon, weights = "stats_C5")

diff_graph <- buildDiffGraph(C1_network, C5_network)
diff_graph <- addCentrality(diff_graph)
diff_graph <- normalizeCentrality(diff_graph)
rank_table <- rankTfs(diff_graph)

library(ggplot2)
ggplot(rank_table, aes(x = rank, y = centrality)) +
  geom_point() +
  ggrepel::geom_text_repel(data = head(rank_table, 10), aes(label = tf)) +
  theme_classic()
```



## 5 Session Info

```
sessionInfo()
#> R version 4.2.0 (2022-04-22)
#> Platform: x86_64-pc-linux-gnu (64-bit)
#> Running under: Ubuntu 18.04.6 LTS
#>
#> Matrix products: default
#> BLAS:    /usr/local/lib/R/lib/libRblas.so
#> LAPACK:  /usr/local/lib/R/lib/libRlapack.so
#>
#> locale:
#> [1] LC_CTYPE=en_US.UTF-8          LC_NUMERIC=C           LC_TIME=C           LC_COLLATE=C          LC_MONETARY=C
#> [11] LC_MEASUREMENT=C          LC_IDENTIFICATION=C
```

## Multome tutorial - MultiAssayExperiment

```
#>
#> attached base packages:
#> [1] stats4      stats       graphics    grDevices   utils       datasets   methods     base
#>
#> other attached packages:
#> [1] ggplot2_3.3.6           msigdbr_7.5.1          epiregulon_1.0.16
#> [8] GenomeInfoDb_1.34.0     IRanges_2.32.0          S4Vectors_0.36.0
#>
#> loaded via a namespace (and not attached):
#> [1] rappidr_0.3.3           tidyverse_1.2.1         bit64_4.0.5
#> [8] data.table_1.14.4        KEGGREST_1.38.0         RCurl_1.98-1.9
#> [15] RSQLite_2.2.18          shadowtext_0.1.2        artificer.mae_1.3.4
#> [22] assertthat_0.2.1       genomitory_2.1.5        viridis_0.6.2
#> [29] dbplyr_2.2.1           Rgraphviz_2.42.0        igraph_1.3.5
#> [36] bookdown_0.29          annotate_1.76.0         libcoin_1.0-9
#> [43] dsdb.plus_1.3.2        cachem_1.0.6            withr_2.5.0
#> [50] metacommons_1.9.0       MultiAssayExperiment_1.24.0 scran_1.26.0
#> [57] ape_5.6-2                crayon_1.5.1            edgeR_3.40.0
#> [64] viper_0.4.5             rlang_1.0.6              lifecycle_1.0.3
#> [71] artificer.base_1.3.19  BiocFileCache_2.6.0      rsvd_1.0.5
#> [78] aplot_0.1.8             Rhdf5lib_1.20.0          zoo_1.8-11
#> [85] bitops_1.0-7            gson_0.0.9              artificer.se_1.3.4
#> [92] blob_1.2.3              DelayedMatrixStats_1.20.0 stringr_1.4.0
#> [99] scales_1.2.1            memoise_2.0.1            GSEABase_1.60.0
#> [106] compiler_4.2.0          tinytex_0.42             dqrng_0.3.0
#> [113] patchwork_1.1.2        ArtifactDB_1.9.5        MASS_7.3-58.1
#> [120] BiocSingular_1.14.0   locfit_1.5-9.6           ggrepel_0.9.1
#> [127] rstudioapi_0.13        bluster_1.8.0            metapod_1.6.0
#> [134] BiocManager_1.30.19   Rcpp_1.0.9               scuttle_1.8.0
#> [141] XML_3.99-0.12          splines_4.2.0            yulab.utils_0.0.5
#> [148] ArchR_1.0.2            ggplotify_0.1.0          xtable_1.8-4
#> [155] ggfunk_0.0.7            ShadowArray_1.7.1        R6_2.5.1
#> [162] clusterProfiler_4.6.0   BiocParallel_1.32.0      BiocNeighbors_1.16.0
#> [169] utf8_1.2.2              lattice_0.20-45          tibble_3.1.8
#> [176] GO.db_3.16.0            survival_3.4-0           limma_3.54.0
#> [183] rhdf5_2.42.0            GenomeInfoDbData_1.2.9  HDF5Array_1.26.0
```