

motifStack guide

Jianhong Ou, Lihua Julie Zhu

October 17, 2016

Contents

1 Introduction

A sequence logo, based on information theory, has been widely used as a graphical representation of sequence conservation (aka motif) in multiple amino acid or nucleic acid sequences. Sequence motif represents conserved characteristics such as DNA binding sites, where transcription factors bind, and catalytic sites in enzymes. Although many tools, such as seqlogo[?], have been developed to create sequence motif and to represent it as individual sequence logo, software tools for depicting the relationship among multiple sequence motifs are still lacking. We developed a flexible and powerful open-source R/Bioconductor package, motifStack, for visualization of the alignment of multiple sequence motifs.

2 Prepare environment

You will need ghostscript: the full path to the executable can be set by the environment variable R_GSCMD. If this is unset, a GhostScript executable will be searched by name on your path. For example, on a Unix, linux or Mac "gs" is used for searching, and on Windows the setting of the environment variable GSC is used, otherwise commands "gswi64c.exe" then "gswin32c.exe" are tried.

Example on Windows: assume that the gswin32c.exe is installed at C:\Program Files\gs\gs9.06\bin, then open R and try:

```
Sys.setenv(R_GSCMD=file.path("C:", "Program Files", "gs",  
                             "gs9.06", "bin", "gswin32c.exe"))
```

3 Examples of using motifStack

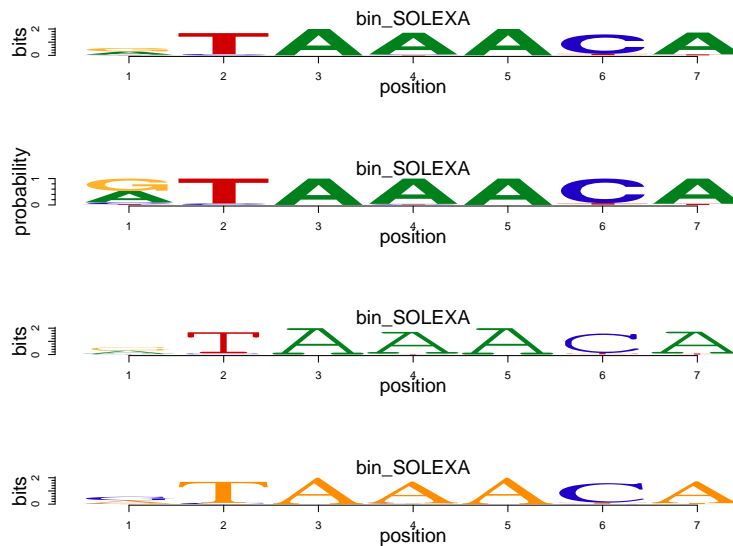


Figure 1: **DNA sequence logo**. Plot a DNA sequence logo with different fonts and colors.

3.1 plot a DNA sequence logo with different fonts and colors

Users can select different fonts and colors to draw the sequence logo (Figure ??).

```
suppressPackageStartupMessages(library(motifStack))
pcm <- read.table(file.path(find.package("motifStack"),
                             "extdata", "bin_SOLEXA.pcm"))

pcm <- pcm[,3:ncol(pcm)]
rownames(pcm) <- c("A", "C", "G", "T")
motif <- new("pcm", mat=as.matrix(pcm), name="bin_SOLEXA")
##pfm object
#motif <- pcm2pfm(pcm)
#motif <- new("pfm", mat=motif, name="bin_SOLEXA")
opar<-par(mfrow=c(4,1))
plot(motif)
#plot the logo with same height
plot(motif, ic.scale=FALSE, ylab="probability")
#try a different font
plot(motif, font="mono,Courier")
#try a different font and a different color group
motif@color <- colorset(colorScheme='basepairing')
plot(motif,font="Times")
par(opar)
```

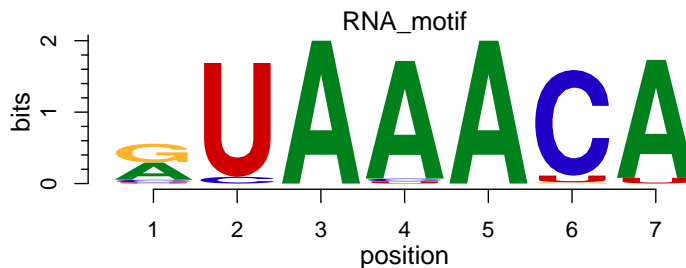


Figure 2: **RNA sequence logo.** Plot an RNA sequence logo

3.2 plot a RNA sequence logo

From DNA sequence logo to RNA sequence logo (Figure ??), you just need to change the rowname of the matrix from "T" to "U".

```
rna <- pcm
rownames(rna)[4] <- "U"
motif <- new("pcm", mat=as.matrix(rna), name="RNA_motif")
plot(motif)
```

3.3 plot an amino acid sequence logo

Given that motifStack allows to use any letters as symbols, it can also be used to draw amino acid sequence logos (Figure ??).

```
library(motifStack)
protein<-read.table(file.path(find.package("motifStack"),"extdata","cap.txt"))
protein<-t(protein[,1:20])
motif<-pcm2pfm(protein)
motif<-new("pfm", mat=motif, name="CAP",
           color=colorset(alphabet="AA",colorScheme="chemistry"))
plot(motif)
```

3.4 plot sequence logo stack

motifStack is designed to show multiple motifs in same canvas. To show the sequence logo stack, the distance of motifs need to be calculated first for example by using `MotIV[?]::motifDistances`, which implemented `STAMP[?]`. After alignment, users can use `plotMotifLogoStack` function to draw sequence logos stack (Figure ??) or use `plotMotifLogoStackWithTree` function to show the distance tree with the sequence logos stack (Figure ??) or use `plotMotifStackWithRadialPhylog` function to plot sequence logo

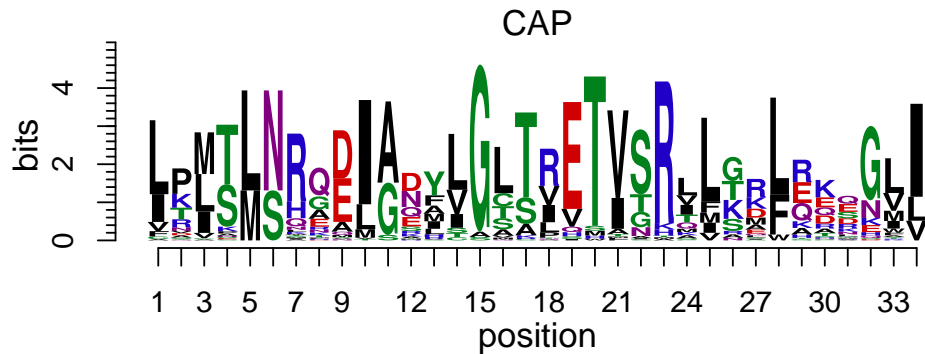


Figure 3: **Amino acid sequence logo.** Plot an sequence logo with any symbols as you want such as amino acid sequence logo

stack in radial style (Figure ??) in the same canvas. There is a shortcut function named as motifStack. Use stack layout to call plotMotifLogoStack, treeview layout to call plotMotifLogoStackWithTree and radialPhylog to call plotMotifStackWithRadialPhylog.

```
library(motifStack)
#####Input#####
pcms<-readPCM(file.path(find.package("motifStack"), "extdata"), "pcm$")
motifs<-lapply(pcms, pcm2pfm)

## plot stacks
motifStack(motifs, layout="stack", ncex=1.0)

## plot stacks with hierarchical tree
motifStack(motifs, layout="tree")

## When the number of motifs is too much to be shown in a vertical stack,
## motifStack can draw them in a radial style.
## random sample from MotifDb
library("MotifDb")
matrix.fly <- query(MotifDb, "Dmelanogaster")
motifs2 <- as.list(matrix.fly)
## use data from FlyFactorSurvey
motifs2 <- motifs2[grepl("Dmelanogaster\\-FlyFactorSurvey\\-",
                        names(motifs2))]

## format the names
names(motifs2) <- gsub("Dmelanogaster_FlyFactorSurvey_", "",
                      gsub("_FBgn\\d+$", "",
                           gsub("[^a-zA-Z0-9]", "_",
```

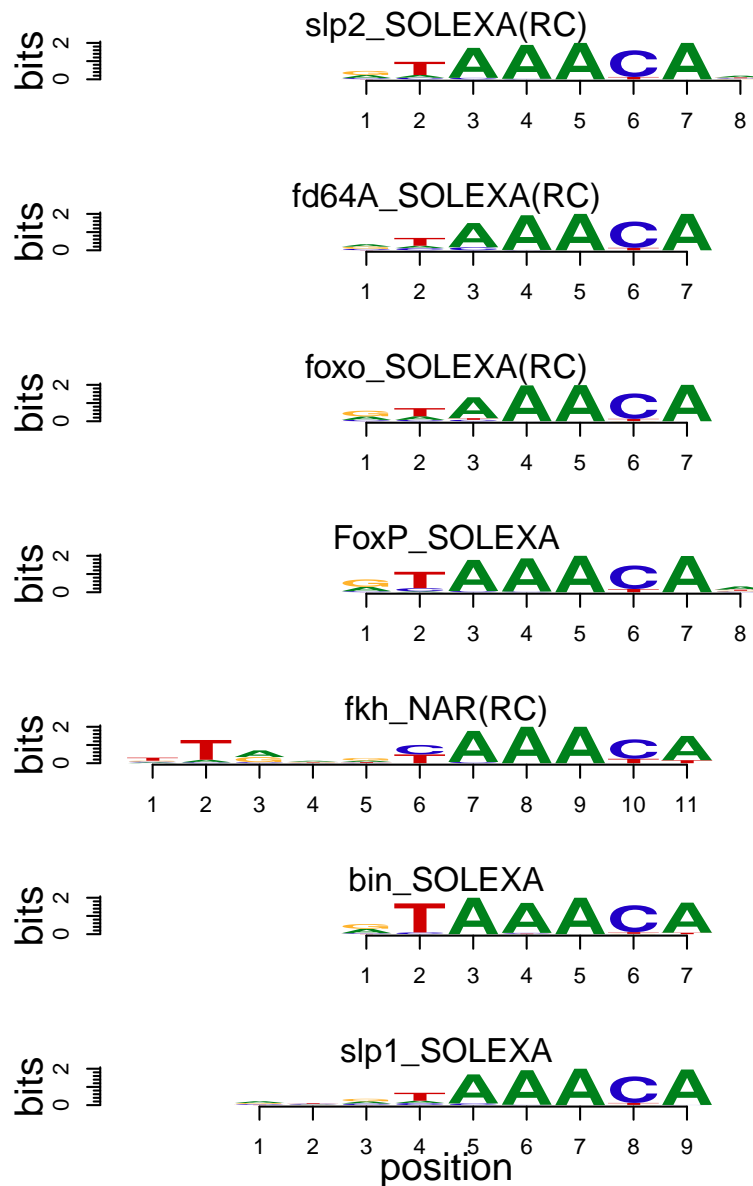


Figure 4: **Sequence logo stack.** Plot motifs with sequence logo stack style.

```

                                gsub("(_\\d+)+$", "", names(motifs2))))))
motifs2 <- motifs2[unique(names(motifs2))]
pfms <- sample(motifs2, 50)
## creat a list of object of pfm
motifs2 <- lapply(names(pfms),
                  function(.ele, pfms){new("pfm",mat=pfms[[.ele]], name=.ele)})
                  ,pfms)
## trim the motifs
motifs2 <- lapply(motifs2, trimMotif, t=0.4)
## setting colors

```

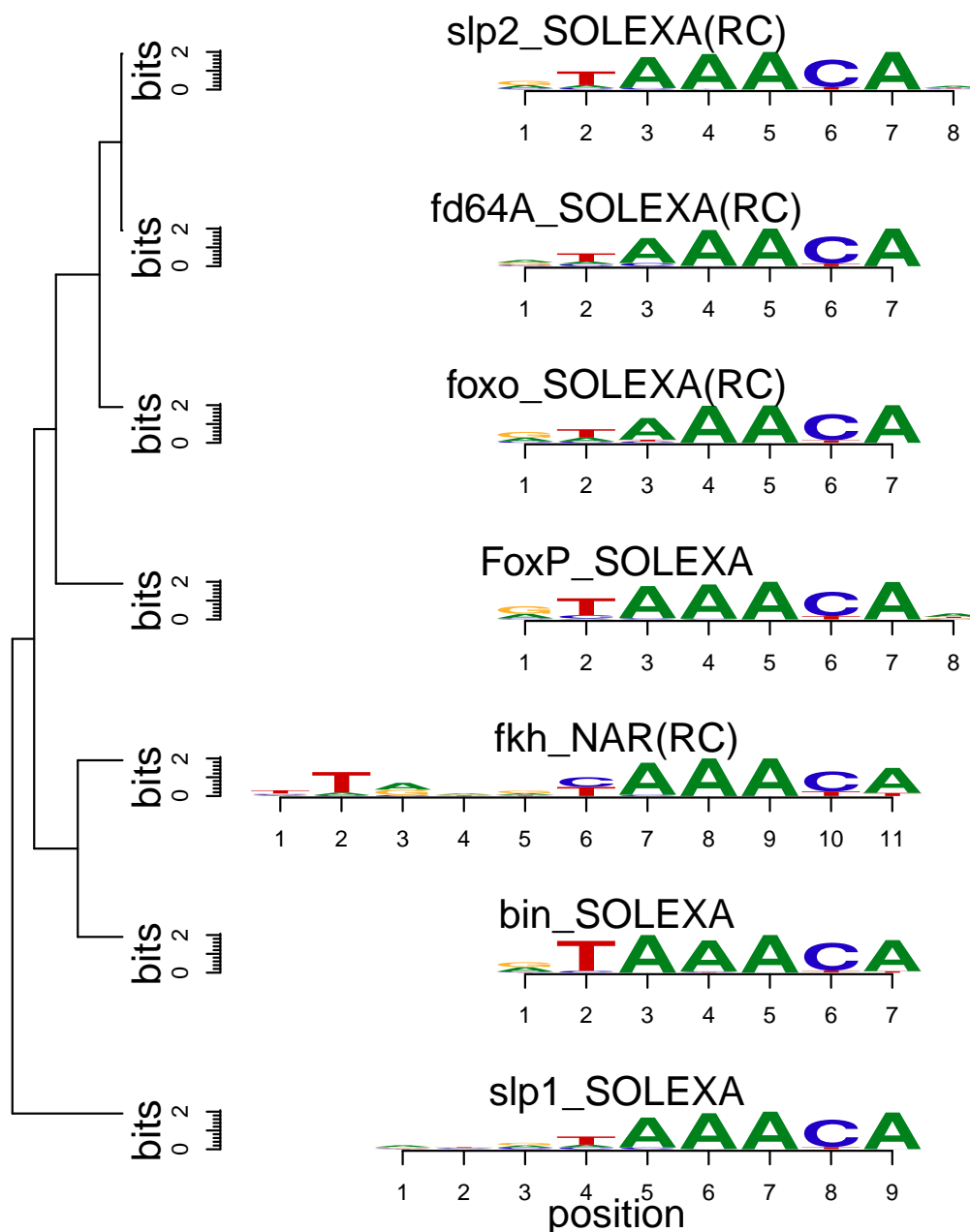


Figure 5: **Treeview layout logo stack.** Sequence logo stack with hierarchical cluster tree.

```
library(RColorBrewer)
color <- brewer.pal(12, "Set3")
## plot logo stack with radial style
motifStack(motifs2, layout="radialPhylog",
  circle=0.3, cleaves = 0.2,
  clabel.leaves = 0.5,
  col.bg=rep(color, each=5), col.bg.alpha=0.3,
  col.leaves=rep(color, each=5),
```

```
col.inner.label.circle=rep(color, each=5),
inner.label.circle.width=0.05,
col.outer.label.circle=rep(color, each=5),
outer.label.circle.width=0.02,
circle.motif=1.2,
angle=350)
```

3.5 plot a sequence logo cloud

We can also plot a sequence logo cloud for DNA sequence logo (Figure ??).

```
## assign groups for motifs
groups <- rep(paste("group",1:5,sep=""), each=10)
names(groups) <- names(pfms)
## assign group colors
group.col <- brewer.pal(5, "Set3")
names(group.col)<-paste("group",1:5,sep="")
## use MotIV to calculate the distances of motifs
jaspar.scores <- MotIV::readDBScores(file.path(find.package("MotIV"),
                                              "extdata",
                                              "jaspar2010_PCC_SWU.scores"))

d <- MotIV::motifDistances(lapply(pfms, pfm2pwm))
hc <- MotIV::motifHclust(d, method="average")
## convert the hclust to phylog object
phylog <- hclust2phylog(hc)
## reorder the pfms by the order of hclust
leaves <- names(phylog$leaves)
pfms <- pfms[leaves]
## create a list of pfm objects
pfms <- lapply(names(pfms), function(.ele, pfms){
  new("pfm",mat=pfms[[".ele"]], name=.ele)}
      ,pfms)
## extract the motif signatures
motifSig <- motifSignature(pfms, phylog, groupDistance=0.01, min.freq=1)
## draw the motifs with a tag-cloud style.
motifCloud(motifSig, scale=c(6, .5),
            layout="rectangles",
            group.col=group.col,
            groups=groups,
            draw.legend=TRUE)
```

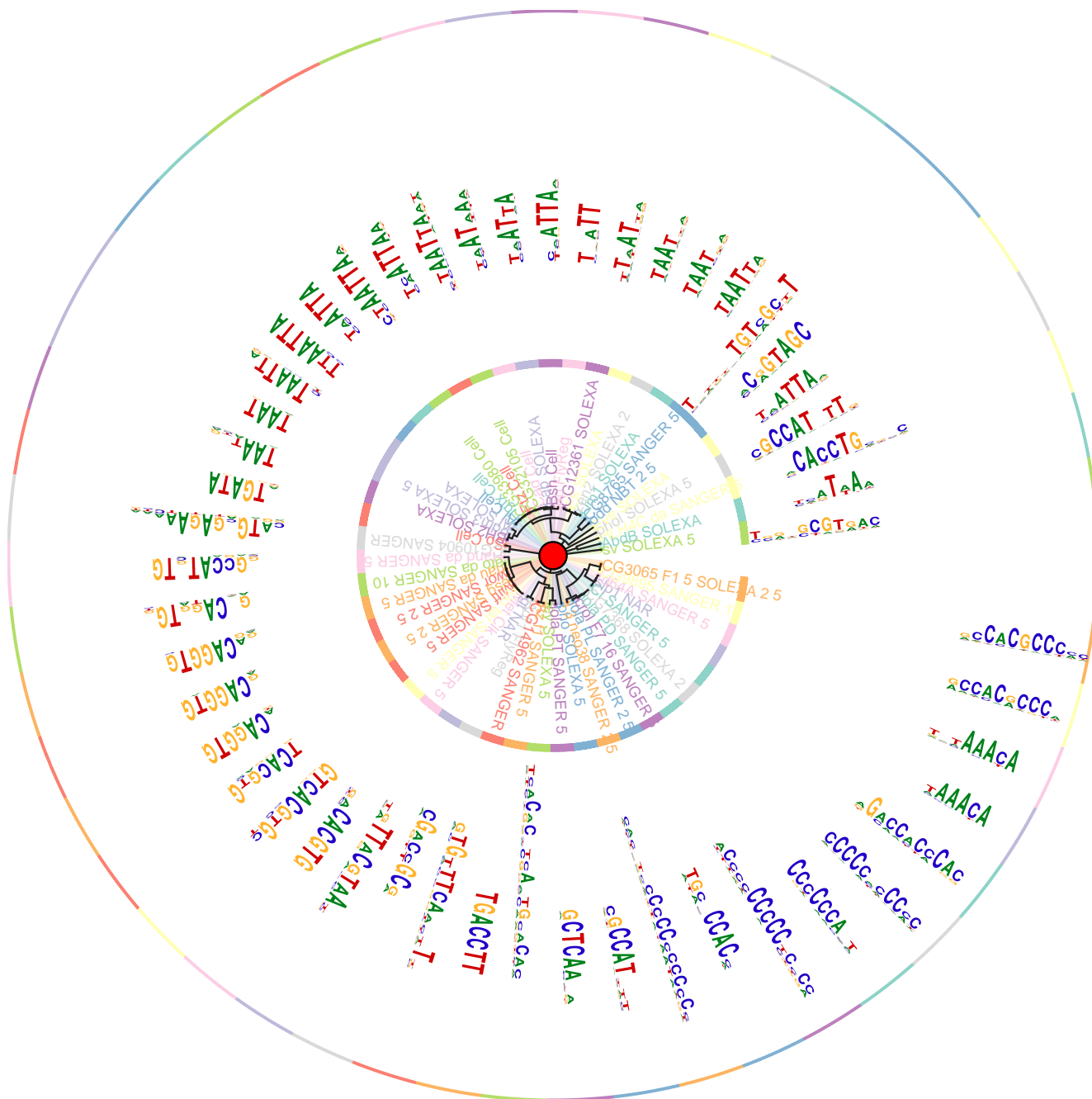


Figure 6: **Sequence logo stack in radial style** Plot motifs in a radial style when the number of motifs is too much to be shown in a vertical stack.

3.6 plot grouped sequence logo

To plot grouped sequence logo, except do motifCloud, we can also plot it with radialPhylog style (Figure ??).

[illegible]

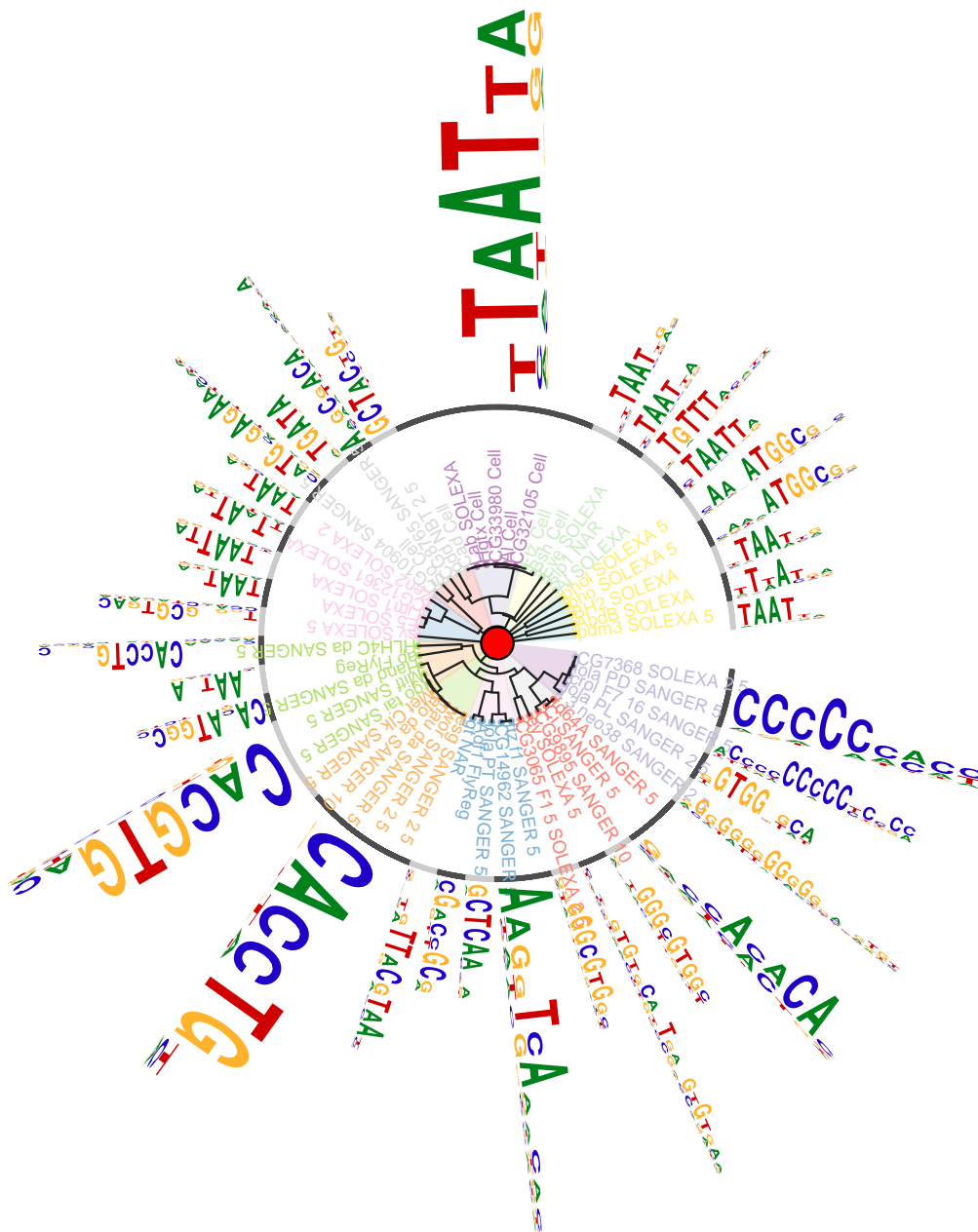


Figure 8: **Grouped sequence logo with radialPhylog style layout.** Like tag-cloud, the sequence logo size is determined by the number of motifs for the signature. The gray-black circle indicates the range of each signature.

```
inner.label.circle.width=0.03,
angle=350, circle.motif=1.2,
motifScale="logarithmic")
```

3.7 motifCircos

We can also plot it with circos style (Figure ??). In circos style, we can plot two group of motifs and with multiple color rings.

```
## plot the logo stack with radial style.
motifCircos(phylog=phylog, pfms=pfms, pfms2=sig,
             col.tree.bg=rep(color, each=5), col.tree.bg.alpha=0.3,
             col.leaves=rep(rev(color), each=5),
             col.inner.label.circle=gpCol,
             inner.label.circle.width=0.03,
             col.outer.label.circle=gpCol,
             outer.label.circle.width=0.03,
             r.rings=c(0.02, 0.03, 0.04),
             col.rings=list(sample(colors(), 50),
                             sample(colors(), 50),
                             sample(colors(), 50)),
             angle=350, motifScale="logarithmic")
```

3.8 motifPiles

We can also plot it with pile style (Figure ??). In pile style, we can plot two group of motifs and with multiple color annotations.

```
## plot the logo stack with radial style.
motifPiles(phylog=phylog, pfms=pfms, pfms2=sig,
            col.tree=rep(color, each=5),
            col.leaves=rep(rev(color), each=5),
            col.pfms2=gpCol,
            r.anno=c(0.02, 0.03, 0.04),
            col.anno=list(sample(colors(), 50),
                           sample(colors(), 50),
                           sample(colors(), 50)),
            motifScale="logarithmic",
            plotIndex=TRUE,
            groupDistance=0.01)
```

4 docker container for motifStack

[Docker](#) allows software to be packaged into containers and the containers can be run any platform as well using a virtual machine called boot2docker. motifStack has its docker image stored in [Docker Hub](#). Users can download the image and run.

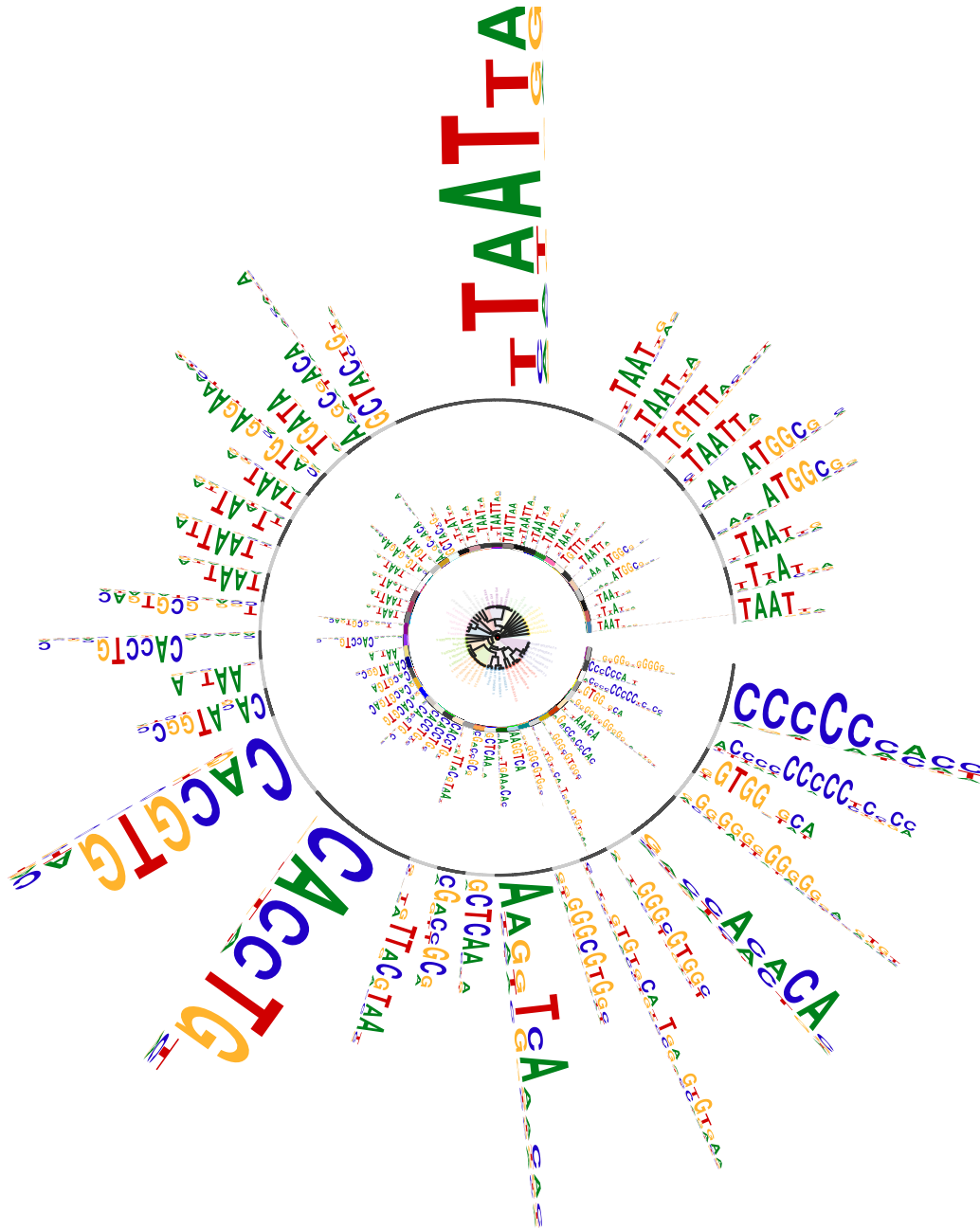


Figure 9: **Grouped sequence logo with circos style layout.** more color sets with more motifs.

```
docker pull jianhong/motifstack_1.13.6
cd ~ ## in windows, please try cd c:\textbackslash Users\textbackslash username
mkdir tmp4motifstack ## this will be the share folder for your host and container.
docker run -ti --rm -v ${PWD}/tmp4motifstack:/volume/data jianhong/motifstack_1.13.6 R
## in R
setwd("/tmp")
library(motifStack)
```

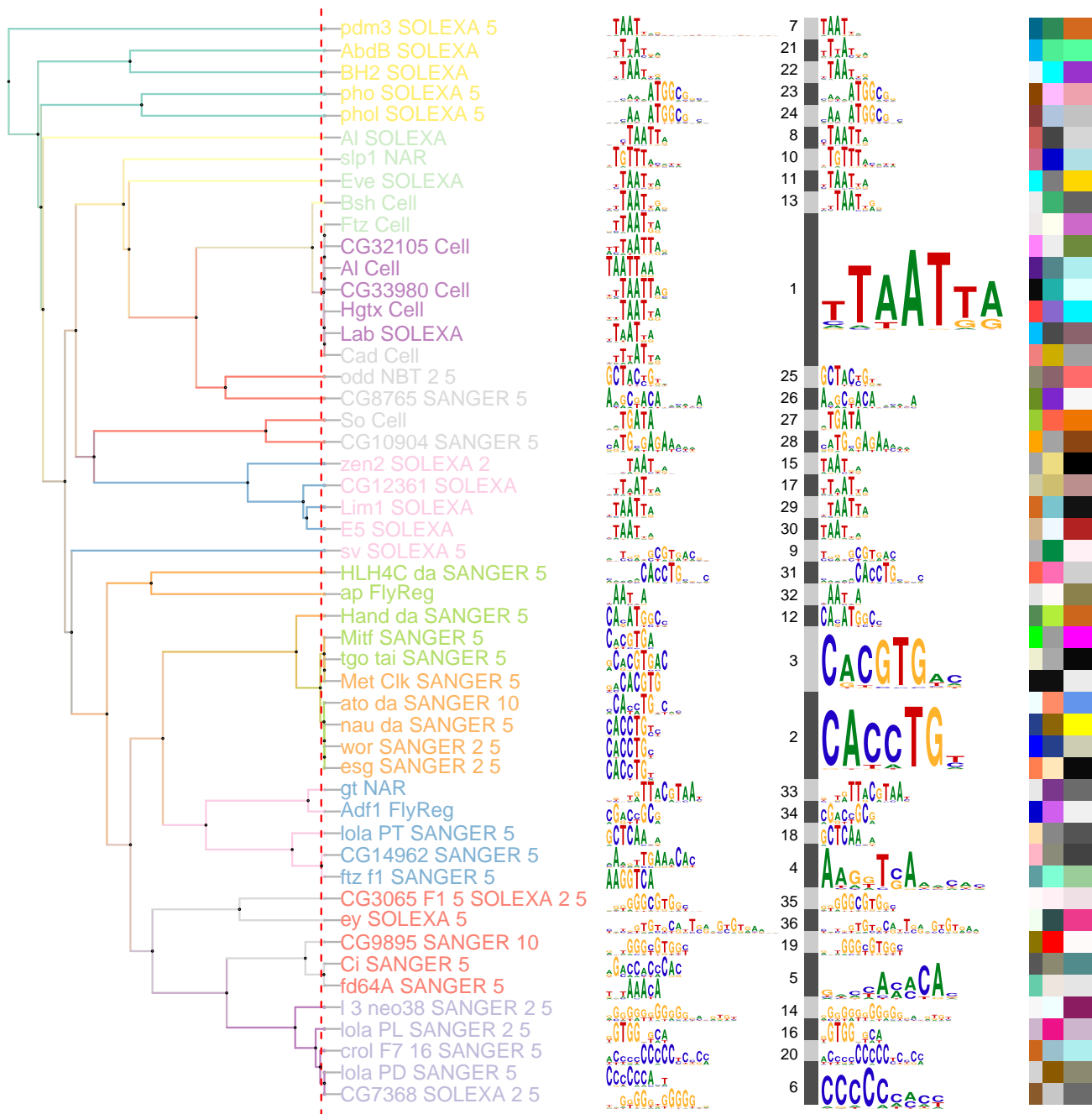


Figure 10: **Grouped sequence logo with piles style layout.** more color sets with more motifs.

```
packageVersion("motifStack")
pcmpath <- "pcmsDatasetFly"
pcms <- readPCM(pcmpath)
pfms <- lapply(pcms, pcm2pfm)
matalign_path <- "/usr/bin/matalign"
neighbor_path <- "/usr/bin/phylip/neighbor"
outpath <- "output"
```

```

system(paste("perl MatAlign2tree.pl --in . --pcmpath", pcmpath, "--out", outpath,
  "--matalign", matalign_path, "--neighbor", neighbor_path, "--tree","UPGMA"))
newickstrUPGMA <- readLines(con=file.path(outpath, "NJ.matalign.distMX.nwk"))
phylog <- newick2phylog(newickstrUPGMA, FALSE)
leaves <- names(phylog$leaves)
motifs <- pfms[leaves]
motifSig <- motifSignature(motifs, phylog, groupDistance=2, min.freq=1, trim=.2)
sig <- signatures(motifSig)
gpCol <- sigColor(motifSig)
leaveNames <- gsub("^Dm_", "", leaves)
pdf("/volume/data/test.pdf", width=8, height=11)
motifPiles(phylog=phylog, DNAmotifAlignment(motifs), sig,
  col.pfms=gpCol, col.pfms.width=.01,
  col.pfms2=gpCol, col.pfms2.width=.01,
  labels.leaves=leaveNames,
  plotIndex=c(FALSE, TRUE), IndexCex=1,
  groupDistance=2, clabel.leaves=1)
dev.off()

```

You will see the test.pdf file in the folder of tmp4motifstack.

5 References

6 Session Info

```

sessionInfo()

## R version 3.3.1 (2016-06-21)
## Platform: x86_64-apple-darwin13.4.0 (64-bit)
## Running under: OS X 10.9.5 (Mavericks)
##
## locale:
## [1] C/en_US.UTF-8/en_US.UTF-8/C/en_US.UTF-8/en_US.UTF-8
##
## attached base packages:
## [1] stats4      parallel    grid        stats       graphics    grDevices   utils       datasets
## [9] methods     base
##
## other attached packages:
## [1] RColorBrewer_1.1-2 MotifDb_1.16.0      motifStack_1.18.0   Biostrings_2.42.0
## [5] XVector_0.14.0      IRanges_2.8.0      S4Vectors_0.12.0    ade4_1.7-4
## [9] MotIV_1.30.0        BiocGenerics_0.20.0 grImport_0.9-0      XML_3.98-1.4
## [13] BiocStyle_2.2.0

```

```
##
## loaded via a namespace (and not attached):
## [1] Rcpp_0.12.7          highr_0.6             plyr_1.8.4
## [4] formatR_1.4          GenomeInfoDb_1.10.0   bitops_1.0-6
## [7] tools_3.3.1          zlibbioc_1.20.0       digest_0.6.10
## [10] evaluate_0.10        tibble_1.2            lattice_0.20-34
## [13] BSgenome_1.42.0      Matrix_1.2-7.1        yaml_2.1.13
## [16] seqLogo_1.40.0       rtracklayer_1.34.0    stringr_1.1.0
## [19] knitr_1.14           Biobase_2.34.0        BiocParallel_1.8.0
## [22] rGADEM_2.22.0        rmarkdown_1.1         magrittr_1.5
## [25] scales_0.4.0         Rsamtools_1.26.0      htmltools_0.3.5
## [28] GenomicRanges_1.26.0 GenomicAlignments_1.10.0 assertthat_0.1
## [31] SummarizedExperiment_1.4.0 colorspace_1.2-7      stringi_1.1.2
## [34] munsell_0.4.3        RCurl_1.95-4.8
```