

dagLogo guide

Jianhong Ou, Lihua Julie Zhu

October 13, 2015

Contents

1	Introduction	1
2	Prepare environment	2
3	Examples of using dagLogo	2
3.1	Step 1, fetch sequences	2
3.2	Step 2, build background model	4
3.3	Step 3, do test	5
3.4	Step 4, graphical representation results	5
4	using dagLogo to analysis Catobolite Activator Protein	12
5	References	15
6	Session Info	16

1 Introduction

A sequence logo has been widely used as a graphical representation of an alignment of multiple amino acid or nucleic acid sequences. There is a package seqlogo[1] implemented in R to draw DNA sequence logos. And another package motifStack[4] was developed for drawing sequence logos for Amino Acid, DNA and RNA sequences. motifStack also has the capability for graphical representation of multiple motifs.

IceLogo[2] is a tool developed in java to visualize significant conserved sequence patterns in an alignment of multiple peptide sequence against background sequences. Compare to webLogo[3], which relying on information theory, iceLogo builds on probability theory. It is reported that iceLogo has a more dynamic nature and is correcter and completer in the analysis of conserved sequence patterns.

However iceLogo can only compare conserved sequences to reference sequences peptide by peptide. As we know, some conserved sequence patterns are not conserved by peptides but by groups such as charge, chemistry, hydrophobicity and etc.

Here we developed a R package:dagLogo based on iceLogo to visualize significant conserved sequence patterns in groups.

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 dagLogo

3.1 Step 1, fetch sequences

You should have interesting peptides position info and the identifiers for fetching sequences via biomaRt.

```

suppressPackageStartupMessages(library(dagLogo))
library(biomaRt)
mart <- useMart("ensembl", "dmelanogaster_gene_ensembl")
dat <- read.csv(system.file("extdata", "dagLogoTestData.csv",
                           package="dagLogo"))
dat <- dat[1:5,] ##subset to speed sample
dat

##   entrez_geneid  NCBI_acc  NCBI_site  molecular_weight
## 1         44149  NP_524708      K328         85400.81Da
## 2         44149  NP_524708       K43         85400.81Da
## 3         44149  NP_524708      K123         85400.81Da
## 4         44149  NP_524708      K409         85400.81Da
## 5         44149  NP_524708      K446         85400.81Da
##
##                peptide peptide..7..7.
## 1          AGIASEAQk*YQA  GIASEAQkYQAKILS
## 2          ALSk*FSDSDVYLPYEK  ASKVALSkFSDSDVYL
## 3    AMQDATAQMALLQFISSGLk*K  QFISSGLkKVAVPST
## 4          CASIAk*DAMSHGLK  GRCASIAkDAMSHGL
## 5  DGISEVFDk*FGGTVLANACGPCIGQWDR  GISEVFDkFGGTVLA

try({
  seq <- fetchSequence(as.character(dat$entrez_geneid),
                       anchorPos=as.character(dat$NCBI_site),
                       mart=mart,
                       upstreamOffset=7,
                       downstreamOffset=7)
  head(seq@peptides)
})

##      [,1] [,2] [,3] [,4] [,5] [,6] [,7] [,8] [,9] [,10] [,11] [,12] [,13]
## [1,] "G"  "I"  "A"  "S"  "E"  "A"  "Q"  "K"  "Y"  "Q"  "A"  "K"  "I"
## [2,] "A"  "S"  "K"  "V"  "A"  "L"  "S"  "K"  "F"  "D"  "S"  "D"  "V"
## [3,] "Q"  "F"  "I"  "S"  "S"  "G"  "L"  "K"  "K"  "V"  "A"  "V"  "P"
## [4,] "G"  "R"  "C"  "A"  "S"  "I"  "A"  "K"  "D"  "A"  "M"  "S"  "H"
## [5,] "G"  "I"  "S"  "E"  "V"  "F"  "D"  "K"  "F"  "G"  "G"  "T"  "V"
##
##      [,14] [,15]
## [1,] "L"   "S"

```

```
## [2,] "Y" "L"
## [3,] "S" "T"
## [4,] "G" "L"
## [5,] "L" "A"
```

Sometimes you may already have the peptides sequences in hand. You will use `formatSequence` function to prepare an object of `dagPeptides` for further testing. To use `formatSequence`, you need prepare the proteome by `prepareProteome` function.

```
dat <- unlist(read.delim(system.file("extdata", "grB.txt",
                                   package="dagLogo"),
               header=F, as.is=TRUE))
head(dat)

##                               V11                               V12
## "GHISVKEPTPSIASDISLPATQELRQLR" "EREMFDKASLKLGLDKAVLQMSGRENATN"
##                               V13                               V14
## "XXXXXXMSDIVVVTDLIAVGLKRGSDLLS" "GQDQEEEEIEDILMDTEELMRAEDTEQLK"
##                               V15                               V16
## "ESYATDNEKMTSTPETLLEEIEAKNRELIA" "VENKERTLKRLLLQDQENSLQDNRTSSDSP"

##prepare proteome from a fasta file
proteome <- prepareProteome(
  fasta=system.file("extdata", "HUMAN.fasta",
                    package="dagLogo"))
##prepare object of dagPeptides
seq <- formatSequence(seq=dat, proteome=proteome,
                      upstreamOffset=14, downstreamOffset=15)
```

3.2 Step 2, build background model

Once you have an object of `dagPeptides` in hand, you can start to build background model for DAG test. The background could be random subsequence of whole proteome or your inputs. If the background was built from whole proteome or proteome without your inputs, an object of `Proteome` is required.

To prepare a proteome, there are two methods, from a fasta file or from UniProt webservice. Last example shows how to prepare proteome from a fasta file. Here we show how to prepare proteome via UniProt webservice.

```
if(interactive()){  
  library(UniProt.ws)  
  taxId(UniProt.ws) <- 9606  
  proteome <- prepareProteome(UniProt.ws=UniProt.ws)  
}
```

Then the proteome can be used for background model building.

```
bg <- buildBackgroundModel(seq, bg="wholeGenome",  
                           proteome=proteome)
```

3.3 Step 3, do test

Test can be done without any change of the symbol pattern or with changes of grouped peptides by such as charge, chemistry, hydrophobicity and etc.

```
t0 <- testDAU(seq, bg)  
t1 <- testDAU(seq, bg, group="classic")  
t2 <- testDAU(seq, bg, group="charge")  
t3 <- testDAU(seq, bg, group="chemistry")  
t4 <- testDAU(seq, bg, group="hydrophobicity")
```

3.4 Step 4, graphical representation results

We can use heatmap (Figure 1) or logo (Figure 2,3,4,5) to show the results.

```
dagHeatmap(t0)
```

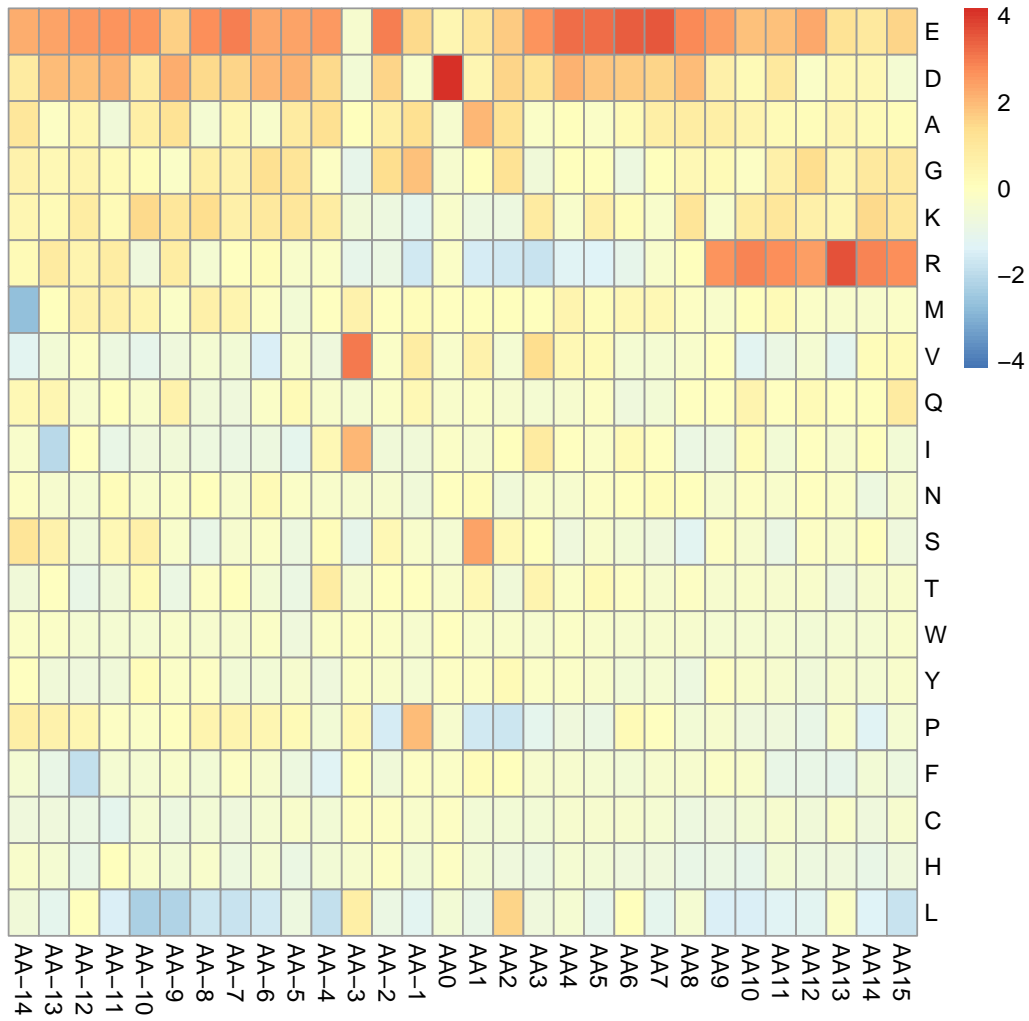


Figure 1: Plot a heatmap to show the results

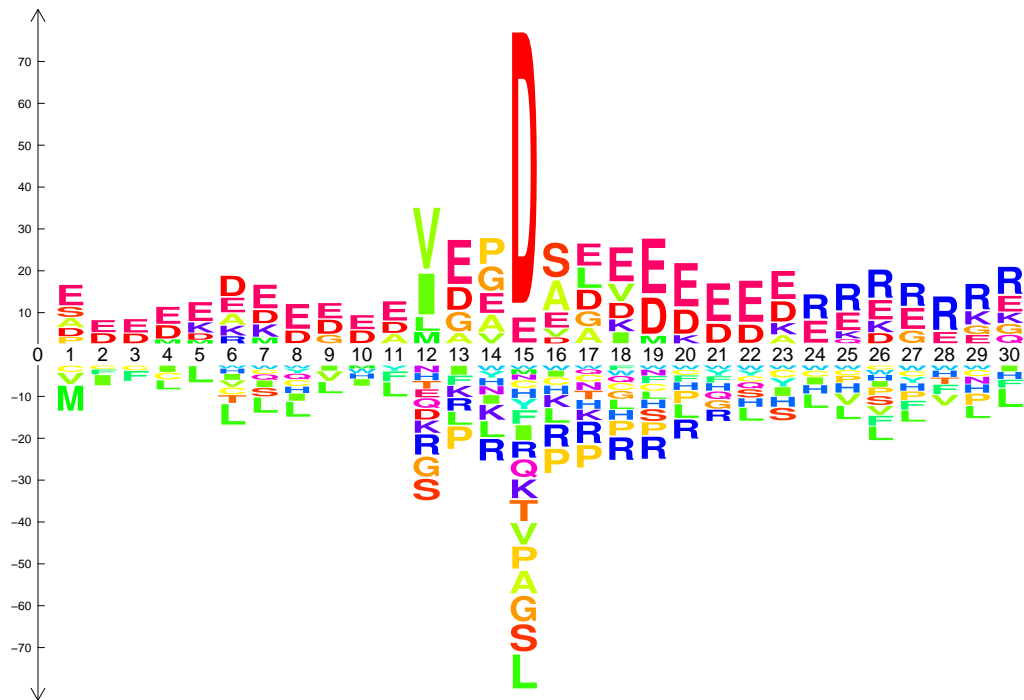


Figure 2: Plot a logo to show the ungrouped results

```
dagLogo(t0)
```

```
dagLogo(t1, namehash=nameHash(t1@group), legend=TRUE)
```

```
dagLogo(t2, namehash=nameHash(t2@group), legend=TRUE)
```

```
dagLogo(t3, namehash=nameHash(t3@group), legend=TRUE)
```

```
dagLogo(t4, namehash=nameHash(t4@group), legend=TRUE)
```

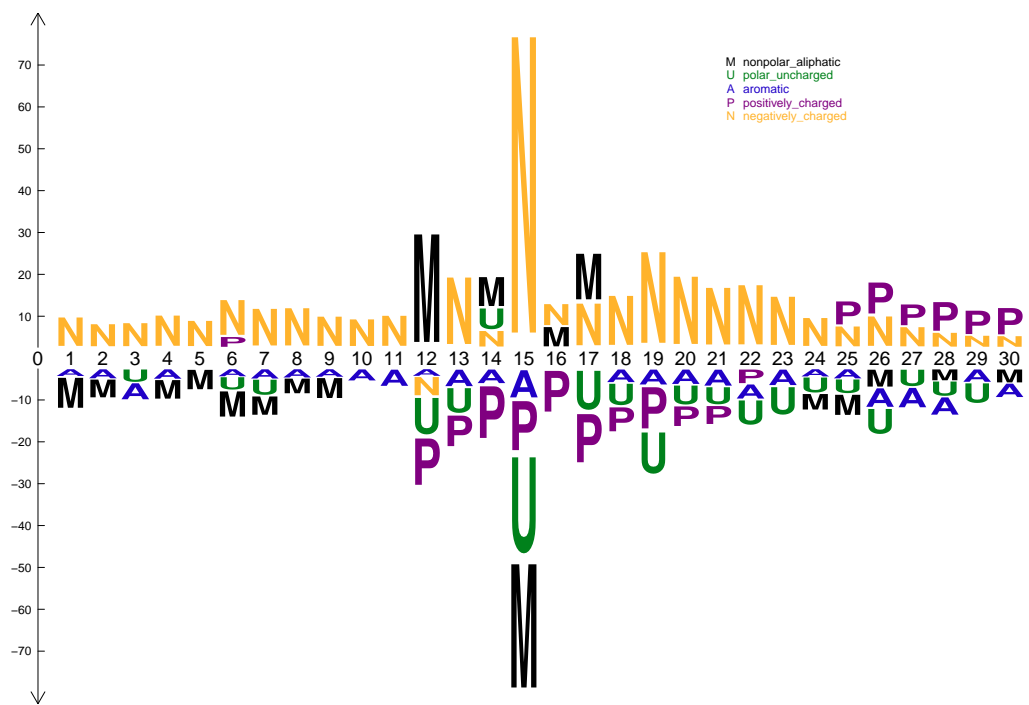


Figure 3: Plot a logo to show the classic grouped results

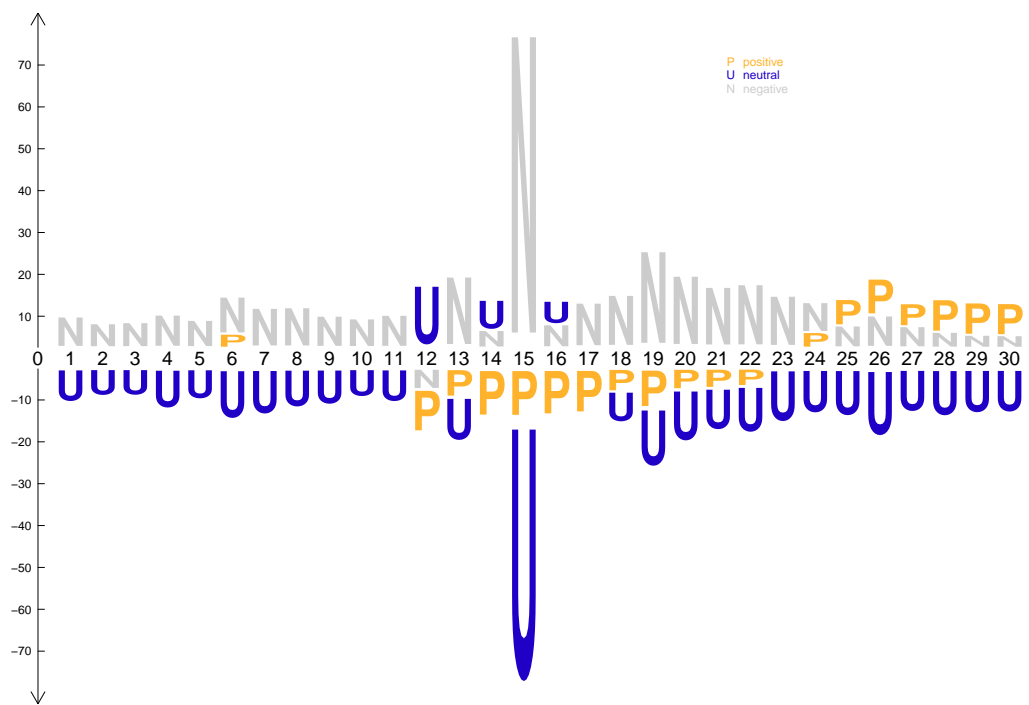


Figure 4: Plot a logo to show the charge grouped results

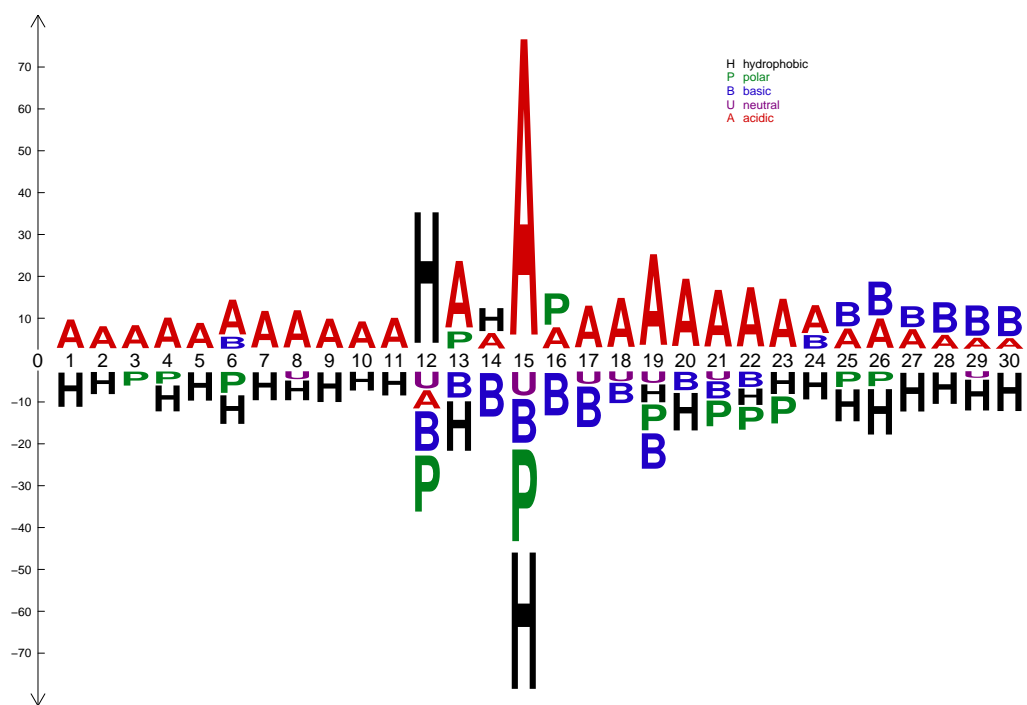


Figure 5: Plot a logo to show the chemistry grouped results

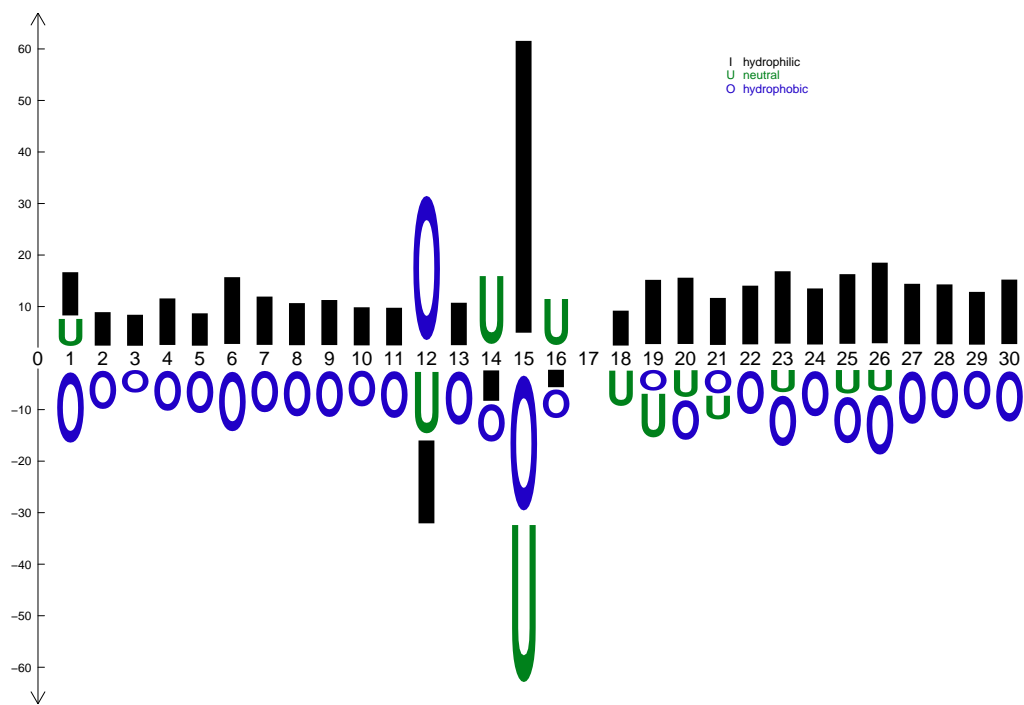


Figure 6: Plot a logo to show the hydrophobicity grouped results

4 using dagLogo to analysis Catobolite Activator Protein

CAP (Catabolite Activator Protein, also known as CRP for cAMP Receptor Protein) is a transcription promoter that binds at more than 100 sites within the E. coli genome.

The motif of the DNA-binding helix-turn-helix motif of the CAP family is drawn by motifStack as Figure 7.

```
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)
```

If we use dagLogo to plot the motif, it will be shown as Figure 8. Residues 7-13 form the first helix, 14-17 the turn and 18-26 the DNA recognition helix. The glycine at position 15 appears to be critical in forming the turn.

```
library(Biostrings)
cap <- as.character(readAAStringSet(
  system.file("extdata", "cap.fasta", package="dagLogo")))
data(ecoli.proteome)
seq <- formatSequence(seq=cap, proteome=ecoli.proteome)
bg <- buildBackgroundModel(seq, bg="wholeGenome",
                          proteome=ecoli.proteome,
                          permutationSize=10L)
t0 <- testDAU(seq, bg)
dagLogo(t0)
```

If the peptides are grouped by chemistry and then plot, it will be shown as Figure 9. Positions 10, 14, 16, 21 and 25 are partially or completely buried and therefore tend to be populated by hydrophobic amino acids, which are very clear if we group the peptides by chemistry.

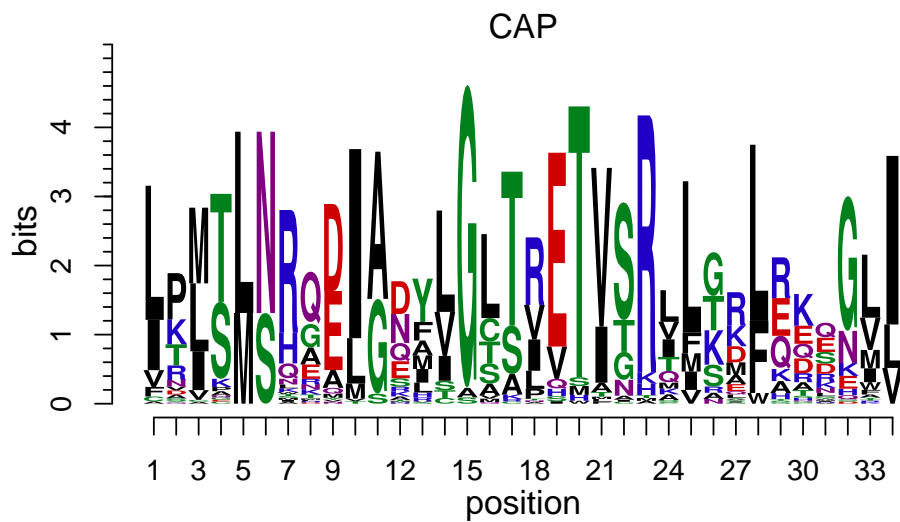


Figure 7: The DNA-binding helix-turn-helix motif of the Catobolite Activator Protein Motif (CAP) family plotted by motifStack

- [3] Gavin E. Crooks, Gary Hon, John-Marc Chandonia, and Steven E. Brenner. Weblogo: A sequence logo generator. *Genome Research*, 14:1188–1190, 2004.
- [4] Jianhong Ou. motifstack: Plot stacked logos for single or multiple dna, rna and amino acid sequence. *R package version 1.5.4*, 2012.

6 Session Info

```

sessionInfo()

## R version 3.2.2 (2015-08-14)
## Platform: x86_64-pc-linux-gnu (64-bit)
## Running under: Ubuntu 14.04.3 LTS
##
## locale:
##  [1] LC_CTYPE=en_US.UTF-8      LC_NUMERIC=C
##  [3] LC_TIME=en_US.UTF-8      LC_COLLATE=C
##  [5] LC_MONETARY=en_US.UTF-8  LC_MESSAGES=en_US.UTF-8
##  [7] LC_PAPER=en_US.UTF-8    LC_NAME=C
##  [9] LC_ADDRESS=C             LC_TELEPHONE=C
## [11] LC_MEASUREMENT=en_US.UTF-8 LC_IDENTIFICATION=C
##
## attached base packages:
##  [1] stats4      parallel  grid      stats      graphics  grDevices  utils
##  [8] datasets  methods  base
##
## other attached packages:
##  [1] knitr_1.11      dagLogo_1.8.0      motifStack_1.14.0
##  [4] Biostrings_2.38.0 XVector_0.10.0     IRanges_2.4.0
##  [7] S4Vectors_0.8.0  ade4_1.7-2         MotIV_1.26.0
## [10] BiocGenerics_0.16.0 grImport_0.9-0     XML_3.98-1.3
## [13] biomaRt_2.26.0    BiocStyle_1.8.0
##
## loaded via a namespace (and not attached):
##  [1] Rcpp_0.12.1      highr_0.5.1

```



```

## [3] RColorBrewer_1.1-2      plyr_1.8.3
## [5] formatR_1.2.1          futile.logger_1.4.1
## [7] GenomeInfoDb_1.6.0     bitops_1.0-6
## [9] futile.options_1.0.0   tools_3.2.2
## [11] zlibbioc_1.16.0        digest_0.6.8
## [13] gtable_0.1.2           RSQLite_1.0.0
## [15] evaluate_0.8           lattice_0.20-33
## [17] BSgenome_1.38.0       DBI_0.3.1
## [19] yaml_2.1.13           seqLogo_1.36.0
## [21] rtracklayer_1.30.0    stringr_1.0.0
## [23] Biobase_2.30.0        AnnotationDbi_1.32.0
## [25] BiocParallel_1.4.0    rGADEM_2.18.0
## [27] rmarkdown_0.8.1      pheatmap_1.0.7
## [29] lambda.r_1.1.7        magrittr_1.5
## [31] scales_0.3.0          GenomicAlignments_1.6.0
## [33] Rsamtools_1.22.0     htmltools_0.2.6
## [35] GenomicRanges_1.22.0 SummarizedExperiment_1.0.0
## [37] colorspace_1.2-6     stringi_0.5-5
## [39] munsell_0.4.2        RCurl_1.95-4.7

```