

canaryGO Analysis

Carson Stacy

2022-09-29

```
# ensure BiocManager is installed, install if not
if (!require("BiocManager", quietly = TRUE))
  install.packages("BiocManager")
BiocManager::install(version = "3.16")

#ensure all required packages are installed
if(!"biomaRt" %in% installed.packages()) BiocManager::install("biomaRt")
if(!"topGO" %in% installed.packages()) BiocManager::install("topGO")
if(!"edgeR" %in% installed.packages()) BiocManager::install("edgeR")
if(!"tidyverse" %in% installed.packages()) install.packages("tidyverse")

library(tidyverse)
library(topGO)
library(biomaRt)
library(here)
library(edgeR)
# library(ggpubr)
```

#TODO: CHECK HOW MANY lncRNA show up in our edgeR outputs. # if they aren't interesting, let's get rid of them and redo our edgeR output. # also probably throw out the tRNAs anyway. # I'll check if BiomaRt works for this.

```
#load in edgeR output to get gene names
# load in edgeR output
MGLpost0 <- readr::read_csv("~/Desktop/canary/contrasts/MGLpost0.csv")
MGpost0 <- readr::read_csv("~/Desktop/canary/contrasts/MGpost0.csv")
MGd140 <- readr::read_csv("~/Desktop/canary/contrasts/MGd140.csv")
MGPpost0 <- readr::read_csv("~/Desktop/canary/contrasts/MGPpost0.csv")
MGd2114 <- readr::read_csv("~/Desktop/canary/contrasts/MGd2114.csv")
Diet0 <- readr::read_csv("~/Desktop/canary/contrasts/Diet0.csv")
MGDietpost0 <- readr::read_csv("~/Documents/GitHub/CanarySeq/contrasts/MGDietpost0.csv")
CtlDietpost0 <- readr::read_csv("~/Documents/GitHub/CanarySeq/contrasts/CtlDietpost0.csv")
CtlLpost0 <- readr::read_csv("~/Documents/GitHub/CanarySeq/contrasts/CtlLpost0.csv")
CtlPpost0 <- readr::read_csv("~/Documents/GitHub/CanarySeq/contrasts/CtlPpost0.csv")

# need to add this contrast for analysis!
#####
# MGDietpost0 <- (((groupMG.lipid.14 + groupMG.lipid.21)/2) - (groupMG.lipid.0)) -
#   (((groupMG.protein.14 + groupMG.protein.21)/2) - (groupMG.protein.0))
#####
# DietMGpost0
```

```

contrasts <- list(
  MGLpost0 = MGLpost0,
  MGpost0 = MGpost0,
  MGd140 = MGd140,
  MGpPost0 = MGpPost0,
  MGd2114 = MGd2114,
  Diet0 = Diet0,
  MGDietpost0 = MGDietpost0,
  # CtlDietpost0 = CtlDietpost0#, # I am checking if this was the contrast with the cool GO term enrichment
  CtlLpost0 = CtlLpost0,
  CtlPpost0 = CtlPpost0
)

```

```

# get biomaRt ensembl dataset for scanaria (has most GO terms)
mart <- if(exists('mart') == TRUE){
  mart
} else {
  useMart(biomart = "ensembl", dataset = "scanaria_gene_ensembl")
}

```

```

# loop through each dataset
for(i in 1:length(contrasts)) {
  # get names of genes for which to get annotations
  transcript_ids <- contrasts[i][[1]][1] #names of all genes in edgeR output

  res <- getBM(attributes = c('ensembl_transcript_id',
                              'ensembl_peptide_id',
                              'ensembl_gene_id',
                              'uniprot_gn_symbol',
                              'wikigene_name',
                              'entrezgene_id',
                              'external_gene_name',
                              'go_id',
                              'description'
                              ),
              filters = 'ensembl_transcript_id',
              values = transcript_ids,
              mart = mart)

  # changing format of the output for subsequent analysis
  # getting GO terms grouped to gene names
  res_min <- res %>%
    mutate(wikigene_name = ifelse(wikigene_name %in% "",
                                   external_gene_name,
                                   wikigene_name))%>% # get gene names
    group_by(ensembl_gene_id) %>% # I think I need to change this...
    mutate(GOterms = paste0(go_id, collapse = ",")) %>%
    slice_head(n=1) %>% # get only one row per gene
    ungroup()

  res_min <- res_min %>%
    mutate(GOterms = gsub("^,+|,$|(|,)+", "\\1", GOterms, perl=TRUE)) %>% # tidy up the extra ;;; in G

```

```

mutate_all(na_if, "")

geneID2GO_combine_names <- res_min[,c(5,7,10)] %>%
  mutate(...1 = ifelse(is.na(res_min$wikigene_name), res_min$external_gene_name, res_min$wikigene_name),
  dplyr::select(-c(wikigene_name, external_gene_name))

geneID_names <- left_join(contrasts[i][[1]][1], geneID2GO_combine_names, by = "...1")
colnames(geneID_names)[1] <- "seq_name"

# build the table of original GO terms with gene IDs
geneID2GO_build_BM <- tibble(
  # canary_annotations$geneID_ncbi,
  geneID_names$seq_name,
  "\t",
  geneID_names$GOterms
)

colnames(geneID2GO_build_BM) <- c("seq_name", "\t", "go_i_ds")

# load in Blast2GO results
Blast2GO <- read_delim("~/Desktop/canary/newGOterms_nomatchDEgenes.tsv",
  delim = "\t", escape_double = FALSE,
  trim_ws = TRUE) %>%
  janitor::clean_names() %>%
  dplyr::select(seq_name, go_i_ds) %>%
  #change formatting of the GO term seperators to ,
  dplyr::mutate(go_i_ds = str_replace_all(go_i_ds, "; ", ", ")) %>%
  dplyr::mutate(go_i_ds = str_replace_all(go_i_ds, "[PFC]:", ""))

combinedGOterms <- left_join(geneID2GO_build_BM, Blast2GO, by = "seq_name")

# now condense the two go_i_ds columns into single column of GO terms.
combinedGOterms <-
  tidyr::unite(combinedGOterms,
    GOterms,
    go_i_ds.x:go_i_ds.y,
    na.rm = TRUE, remove = FALSE,
    sep = ",")
geneID2GO_build_BM <- tibble(
  # canary_annotations$geneID_ncbi,
  seq_name = combinedGOterms$seq_name,
  "\t",
  GOterms = combinedGOterms$GOterms
)

# write out the geneID2GO_BM.map file
write.table(geneID2GO_build_BM, file=here("geneID2GO_BM.map"),
  quote = F, row.names = F, col.names = F)
save(geneID2GO_build_BM,
  file = paste( "~/Documents/GitHub/CanarySeq/geneAnnotations/geneID2GO_", names(contrasts)[i], ".Rsave"))

# read in with readMappings

```

```

canary_geneID2GO_BM <- topGO::readMappings(file = here("geneid2GO_BM.map"))

# save these mappings
save(canary_geneID2GO_BM, file = paste( "~/Documents/GitHub/CanarySeq/geneAnnotations/canary_geneID2GO_BM.RData",
# "~/Desktop/canary/canaryGeneID2GOMappings_BM.RData")

#load in canary geneID2GO map
load(file = paste( "~/Documents/GitHub/CanarySeq/geneAnnotations/canary_geneID2GO_BM_", names(contrasts

# Numeric ID values (in edgeR output) of genes whose names match Gene Ontology Estimates with edgeR o
matches <- which(toupper(names(canary_geneID2GO_BM)) %in% contrasts[i][[1]][[1]])
# these names came from `edgeR diff expression.R` file

# # Numeric ID values (in eggNOG output) of genes whose names match Gene Ontology Estimates with edge
# matches_rnaIndex <- which(
#   contrasts[i][[1]][[1]] %in% toupper(names(canary_geneID2GO_BM))
# )

#genes that are in both.
# used unique because there are some duplicates in output. (per isoform)
genesMatched <- unique(toupper(names(canary_geneID2GO_BM))[matches])

# loop through different lfc cutoffs I want to compare (and up v down)
# for (j in c(1.1,1.2,1.5)) {
for (j in c(1.2)) {
# for (j in c(1.1)) {
# for (j in c(1.5)) {
# for (j in c(2)) {
# for (j in c(1.01,1.1,1.2,1.5)) {
  DEgenes_up <- filter(as.data.frame(contrasts[[i]]), logFC > log2(j)) %>% pull(`...1`)
  DEgenes_down <- filter(contrasts[[i]], logFC < -log2(j)) %>% pull(`...1`)

# DEgenes_up <- filter(as.data.frame(contrasts[[i]]), logFC > j) %>% pull(`...1`)
# DEgenes_down <- filter(contrasts[[i]], logFC < -j) %>% pull(`...1`)

# genes names with GO terms and DE.
DEgeneswithGO_up <- which(DEgenes_up %in% genesMatched)
DEgeneswithGO_down <- which(DEgenes_down %in% genesMatched)

# select only genes that are DE and have GO terms
DEgenesGO_up <- DEgenes_up[DEgeneswithGO_up]
DEgenesGO_down <- DEgenes_down[DEgeneswithGO_down]

#first, let's subset our canary geneIDtoGO to just these genes.
canary_geneID2GO_BM_DE_up <-
  canary_geneID2GO_BM[which(toupper(names(canary_geneID2GO_BM)) %in% DEgenesGO_up)]

names(canary_geneID2GO_BM_DE_up)[!(DEgenesGO_up %in% names(canary_geneID2GO_BM_DE_up))]

canary_geneID2GO_BM_DE_down <-
  canary_geneID2GO_BM[which(toupper(names(canary_geneID2GO_BM)) %in% DEgenesGO_down)]

```

```

names(canary_geneID2GO_BM_DE_down)[!(DEgenesGO_down %in% names(canary_geneID2GO_BM_DE_down))]

#Let me also get my full edgeR gene list.
canary_geneID2GO_BM_edgeR <- canary_geneID2GO_BM[which(toupper(names(canary_geneID2GO_BM)) %in% gene

# weird formatting creating the geneList for topGO. this code seemed to work
all_genes_up <- factor(as.integer(toupper(contrasts[i][[1]][[1]]) %in% DEgenesGO_up)) # create 0 1
names(all_genes_up) <- toupper(contrasts[i][[1]][[1]]) # add gene names

all_genes_down <- factor(as.integer(toupper(contrasts[i][[1]][[1]]) %in% DEgenesGO_down)) # create
names(all_genes_down) <- toupper(contrasts[i][[1]][[1]]) # add gene names

G0data_up=new('topGOdata',
  ontology='BP', # 'BP' = Biological Process. Also tried 'ALL' but doesn't work.
  allGenes = all_genes_up, # all 10188 genes that are in MGpost0 data from edgeR. factor=1 for
  annot = annFUN.gene2GO, # annotation function for custom annotation
  # gene2GO = canary_geneID2GO_BM # annotation db from eggNOG output.
# I am thinking about switching this for one filtered to just edgeR genes:
  gene2GO = canary_geneID2GO_BM_edgeR # annotation db from BiomaRt.
  # gene2GO = canary_geneID2GO_BM_DE#,
  # geneSel = DE_genes
)

G0data_down=new('topGOdata',
  ontology='BP', # 'BP' = Biological Process. Also tried 'ALL' but doesn't work.
  allGenes = all_genes_down, # all 10188 genes that are in MGpost0 data from edgeR. factor=1 for
  annot = annFUN.gene2GO, # annotation function for custom annotation
  # gene2GO = canary_geneID2GO_BM # annotation db from eggNOG output.
# I am thinking about switching this for one filtered to just edgeR genes:
  gene2GO = canary_geneID2GO_BM_edgeR # annotation db from BiomaRt.
  # gene2GO = canary_geneID2GO_BM_DE#,
  # geneSel = DE_genes
)

#analysis
# define test using the weight01 algorithm (default) with fisher
weight_fisher_result_up=runTest(G0data_up, algorithm='weight01', statistic='fisher')
weight_fisher_result_down=runTest(G0data_down, algorithm='weight01', statistic='fisher')

# generate a table of results: we can use the GenTable function to generate a summary table with th
allGO_up=usedGO(G0data_up)
all_res_up=GenTable(G0data_up, weightFisher=weight_fisher_result_up, orderBy='weightFisher', topNod

allGO_down=usedGO(G0data_down)
all_res_down=GenTable(G0data_down, weightFisher=weight_fisher_result_down, orderBy='weightFisher',

#performing BH correction on our p values
p.adj_up=round(p.adjust(all_res_up$weightFisher,method="BH"),digits = 4)
p.adj_down=round(p.adjust(all_res_down$weightFisher,method="BH"),digits = 4)

```

```

# create the file with all the statistics from GO analysis
all_res_final_up=cbind(all_res_up,p.adj_up)
all_res_final_up=all_res_final_up[order(all_res_final_up$p.adj),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

all_res_final_down=cbind(all_res_down,p.adj_down)
all_res_final_down=all_res_final_down[order(all_res_final_down$p.adj),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

#get list of significant GO before multiple testing correction
results.table.p_up= all_res_final_up[which(as.numeric(all_res_final_up$weightFisher)<=0.1),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

results.table.p_down= all_res_final_down[which(as.numeric(all_res_final_down$weightFisher)<=0.1),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

#get list of significant GO after multiple testing correction
results.table.bh_up=all_res_final_up[which(all_res_final_up$p.adj<=0.1),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

results.table.bh_down=all_res_final_down[which(all_res_final_down$p.adj<=0.1),] %>%
  mutate(across(everything(), gsub, pattern = ",", replacement = ";"))

#save first top 50 ontologies sorted by adjusted pvalues
write.table(all_res_final_up[1:50,],
  file = paste(
    "~/Documents/GitHub/CanarySeq/geneAnnotations/geneID2GO_",
    names(contrasts)[i], "log2of", as.character(j), "_up.csv", sep= ""),
  # "~/Desktop/canary/summary_topGO_analysis.csv",
  sep=";",
  quote=FALSE,
  row.names=FALSE)

write.table(all_res_final_down[1:50,],
  file = paste(
    "~/Documents/GitHub/CanarySeq/geneAnnotations/geneID2GO_",
    names(contrasts)[i], "log2of", as.character(j), "_down.csv", sep= ""),
  # "~/Desktop/canary/summary_topGO_analysis.csv",
  sep=";",
  quote=FALSE,
  row.names=FALSE)

# PLOT the GO hierarchy plot: the enriched GO terms are colored in yellow/red according to significance
# pdf(file='~/Desktop/canary/topGOPlot_fullnames.pdf', height=12, width=12, paper='special', points
# showSigOfNodes(GOdata_down, score(weight_fisher_result_up), useInfo = "none", sigForAll=FALSE, fi
# dev.off()
}

```

```

# visualize results at final j value
ntop <- 25

ggdata_up <- all_res_up[1:ntop,] %>%
  mutate(enrichment = paste(names(contrasts[i]), "_UP"))
ggdata_up$Term <- factor(ggdata_up$Term, levels = rev(ggdata_up$Term)) # fixes order

ggdata_down <- all_res_down[1:ntop,] %>%
  mutate(enrichment = paste(names(contrasts[i]), "_DOWN"))
ggdata_down$Term <- factor(ggdata_down$Term, levels = rev(ggdata_down$Term))

ggdata <- bind_rows(ggdata_up, ggdata_down) %>%
  mutate(`% overrepresented` = (Significant-Expected)/Annotated)

gg <- ggplot(ggdata,
  aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = `% overrepresented`))

  expand_limits(y = 1) +
  geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
    linetype = c("dotted", "longdash", "solid", "dotted", "longdash", "solid"),
    colour = "black",
    size = c(0.5, 1.5, 3, 0.5, 1.5, 3)) +
  geom_point(shape = 21) +
  scale_size(range = c(2.5, 12.5)) +
  scale_fill_continuous(low = 'royalblue', high = 'red4') +

  xlab('') +
  ylab('Enrichment score') +
  labs(
    title = paste('GO Biological processes upregulated in contrast: ', names(contrasts[i])),
    subtitle = paste('Top ', ntop, ' terms ordered by weighted Fisher p-value'),
    caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +

  theme_bw(base_size = 24) +
  theme(
    legend.position = 'right',
    legend.background = element_rect(),
    plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust = 0.5),
    plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
    plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),

    axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
    axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
    axis.title = element_text(size = 12, face = 'bold'),
    axis.title.x = element_text(size = 12, face = 'bold'),
    axis.title.y = element_text(size = 12, face = 'bold'),
    axis.line = element_line(colour = 'black'),

    #Legend
    legend.key = element_blank(), # removes the border
    legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend
    legend.text = element_text(size = 14, face = "bold"), # Text size

```

```

    title = element_text(size = 14, face = "bold")) +

    facet_grid(rows = vars(enrichment), scales = "free") +
    coord_flip()

print(gg)

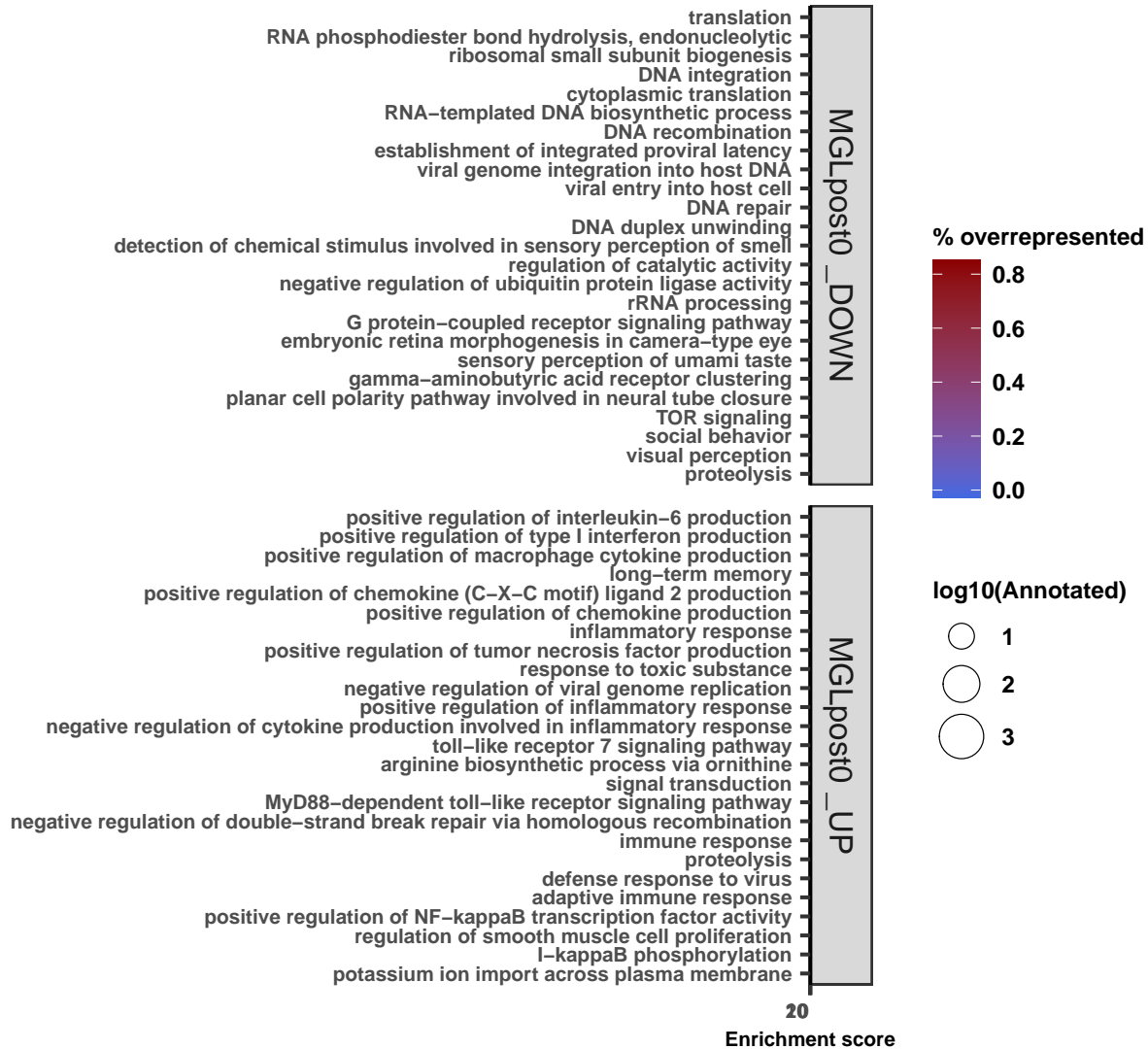
ggsave(gg,
       device = 'png',
       filename = paste("~/Documents/GitHub/CanarySeq/GOanalysis/GO_VIZ_", names(contrasts)[i], ".png", sep = ""))
)

# save GOdata_up and down for i^th contrast at last j cutoff threshold
save(GOdata_up, GOdata_down, DEgenes_up, DEgenes_down,
     file = paste("~/Documents/GitHub/CanarySeq/GOanalysis/GOdataUPDN_", names(contrasts)[i], ".RData", sep = ""))
)
}

```


GO Biological processes upregulated in contrast: MGLposi

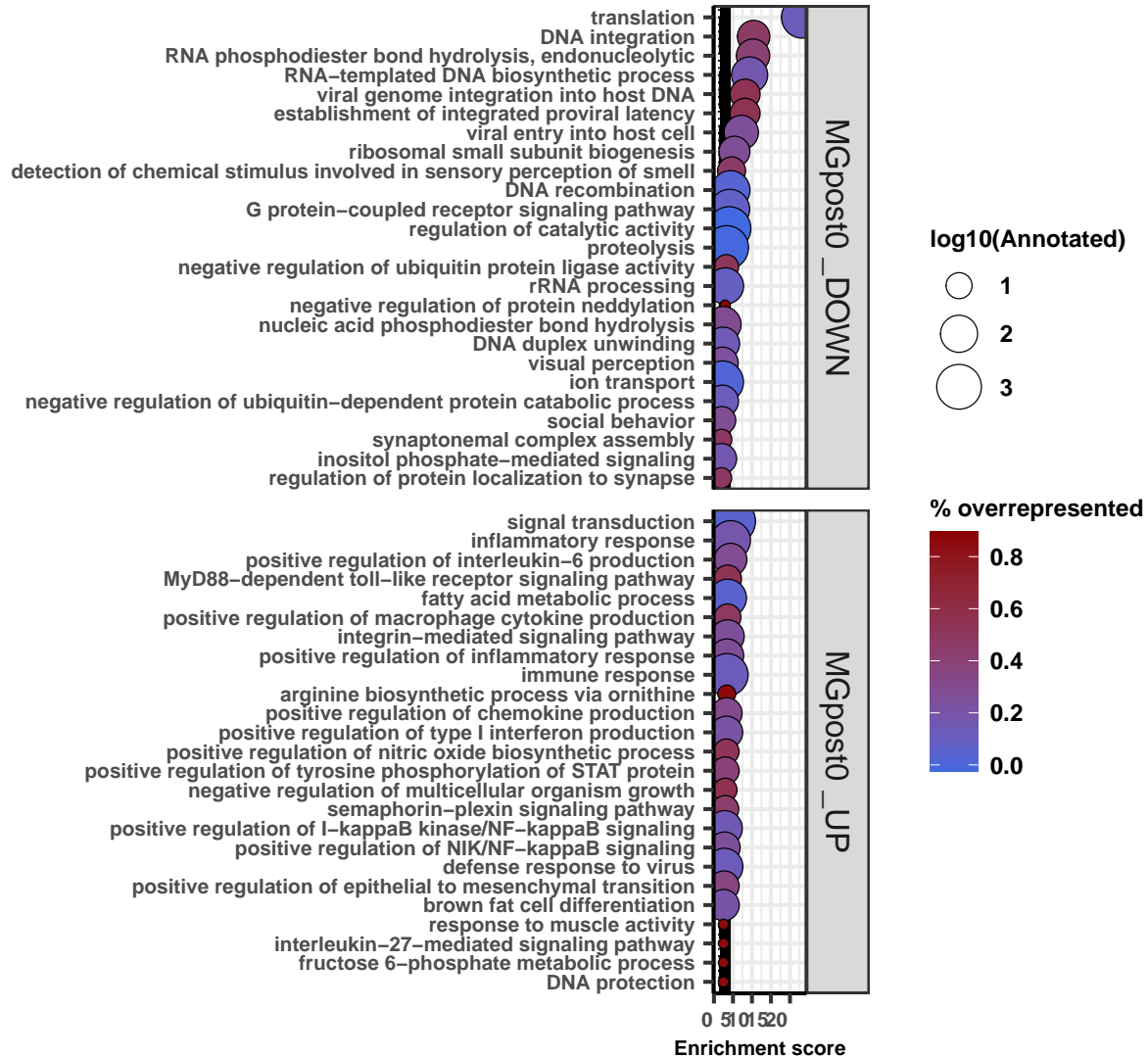
Top 25 terms ordered by weight



Cut-off lines drawn at equivalents of $p=0.05$, $p=0.01$, $p=0.001$

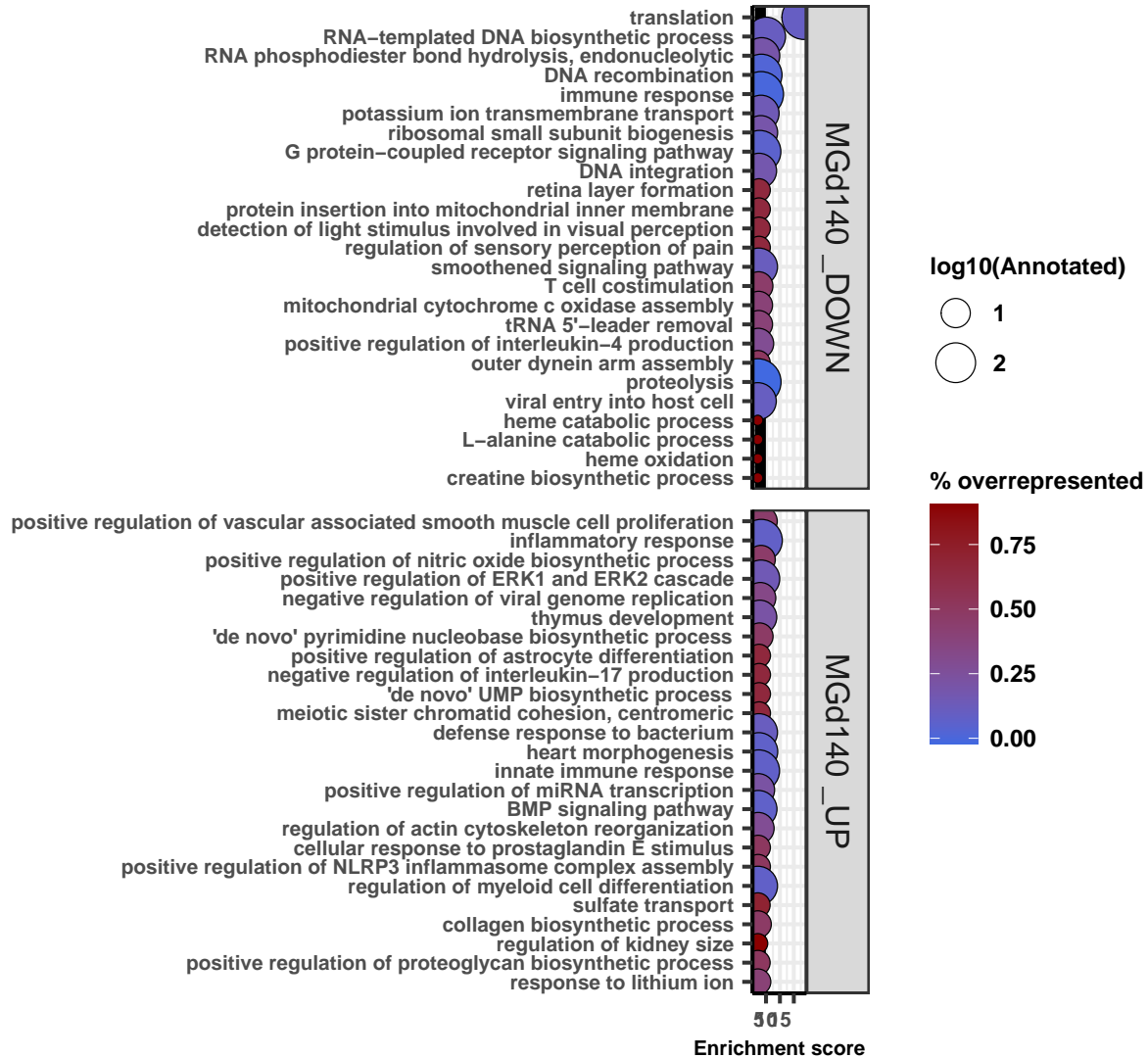
GO Biological processes upregulated in contrast: MGpost0

Top 25 terms ordered by weighted Fisher



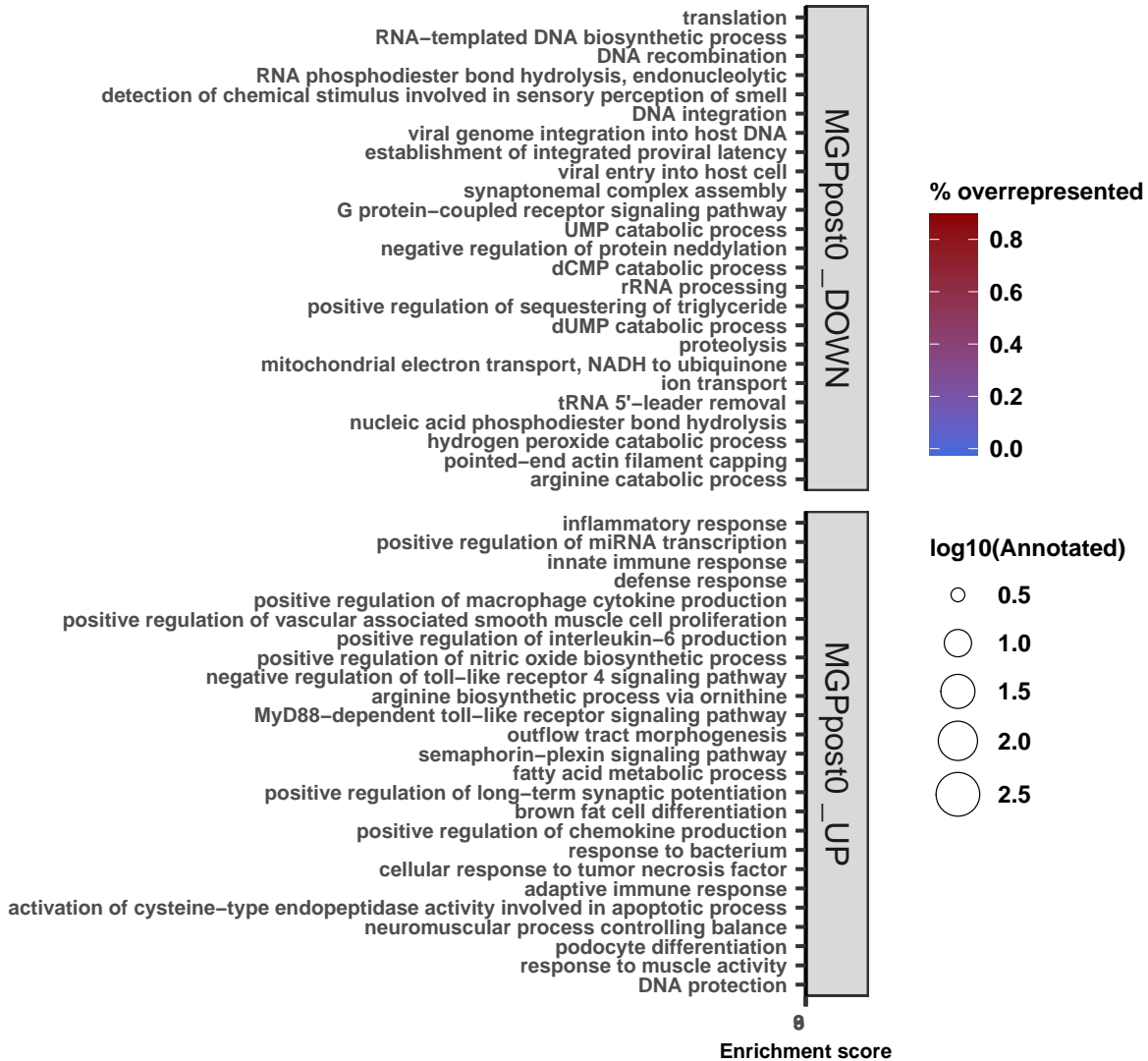
GO Biological processes upregulated in contrast: MGd140

Top 25 terms ordered by weighted Fis



GO Biological processes upregulated in contrast: MGPpost

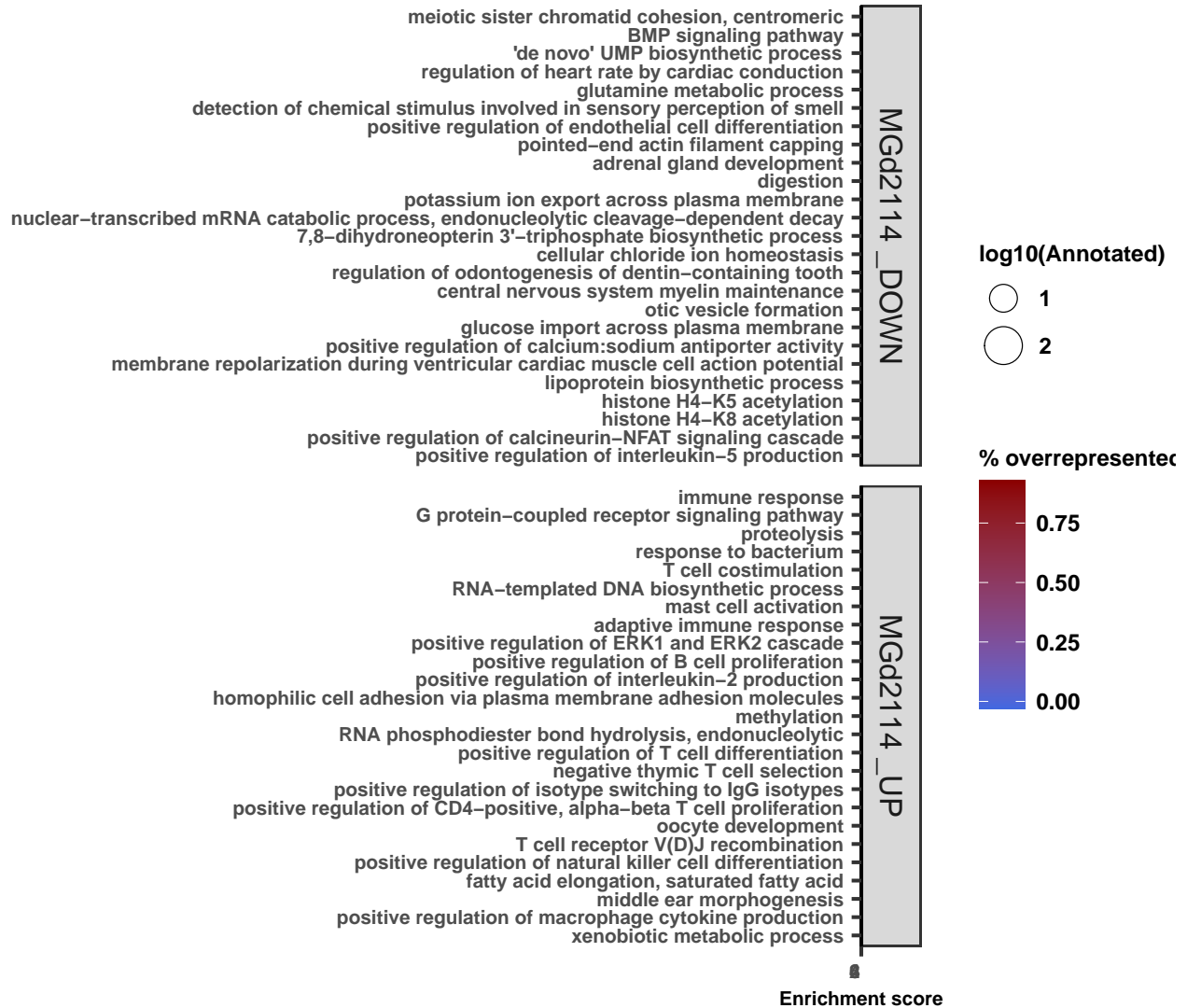
Top 25 terms ordered by weight



Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001

GO Biological processes upregulated in contrast: MGd:

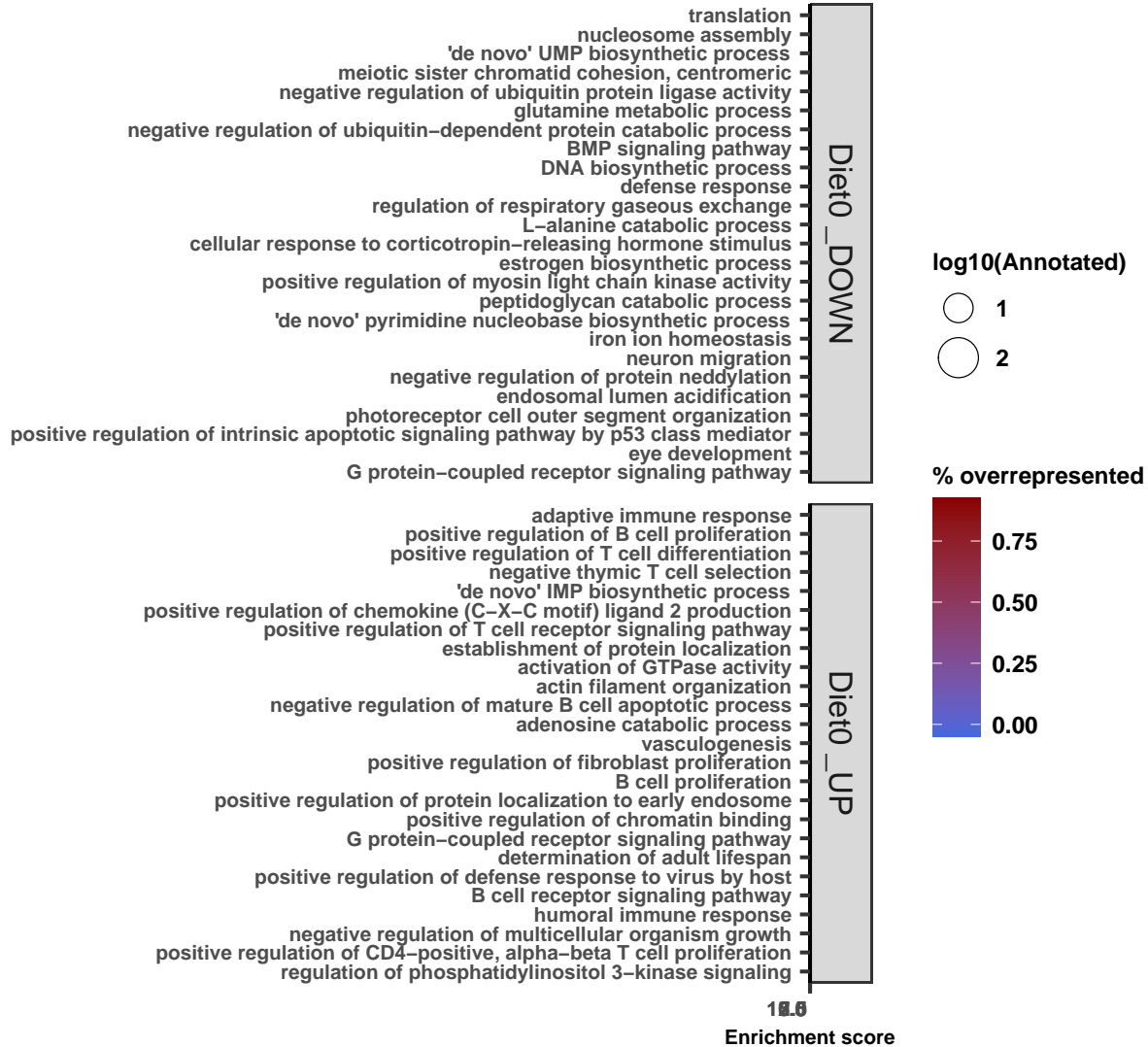
Top 25 terms ordered by wei



Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001

GO Biological processes upregulated in contrast: Diet0

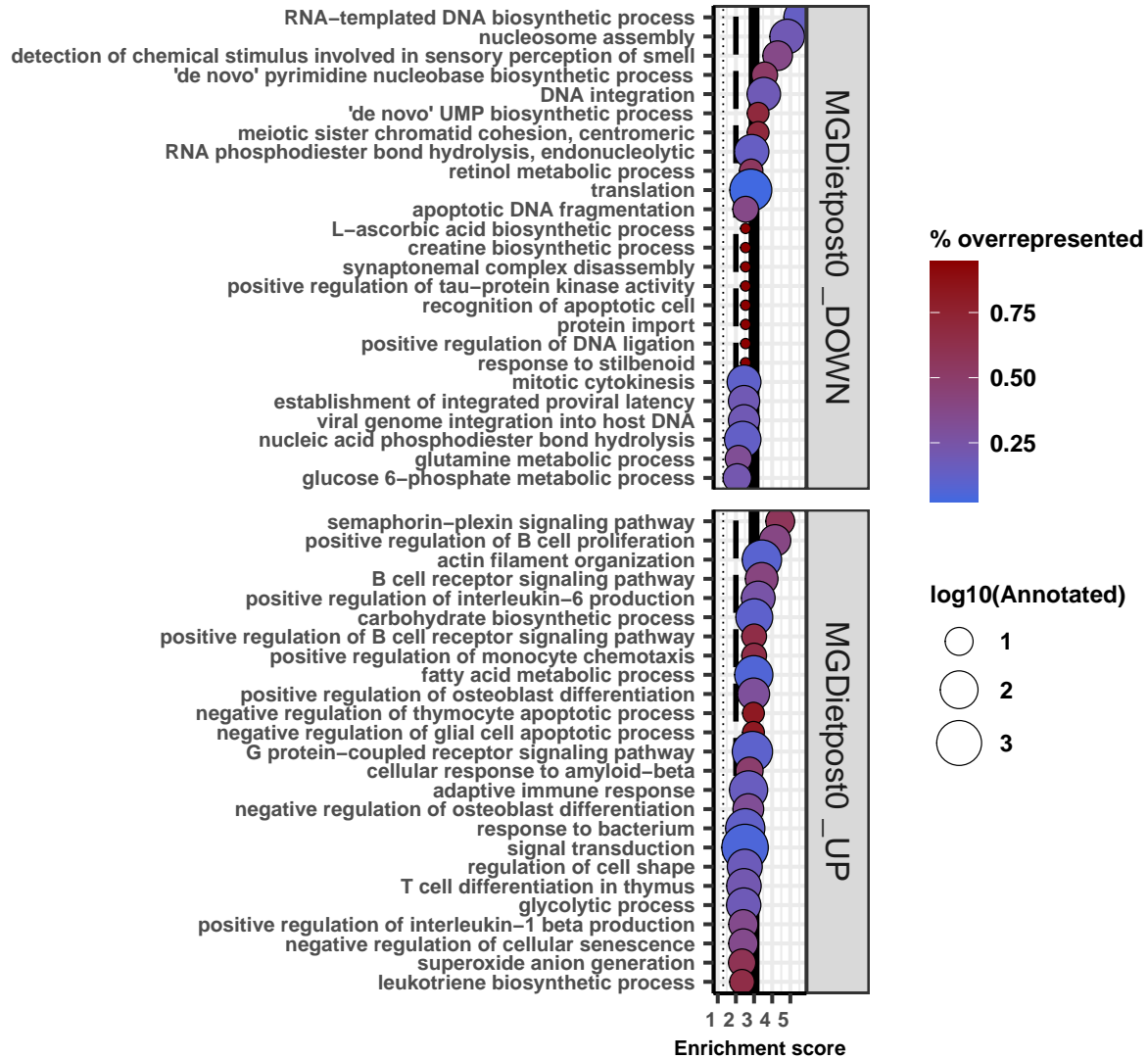
Top 25 terms ordered by weight



Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001

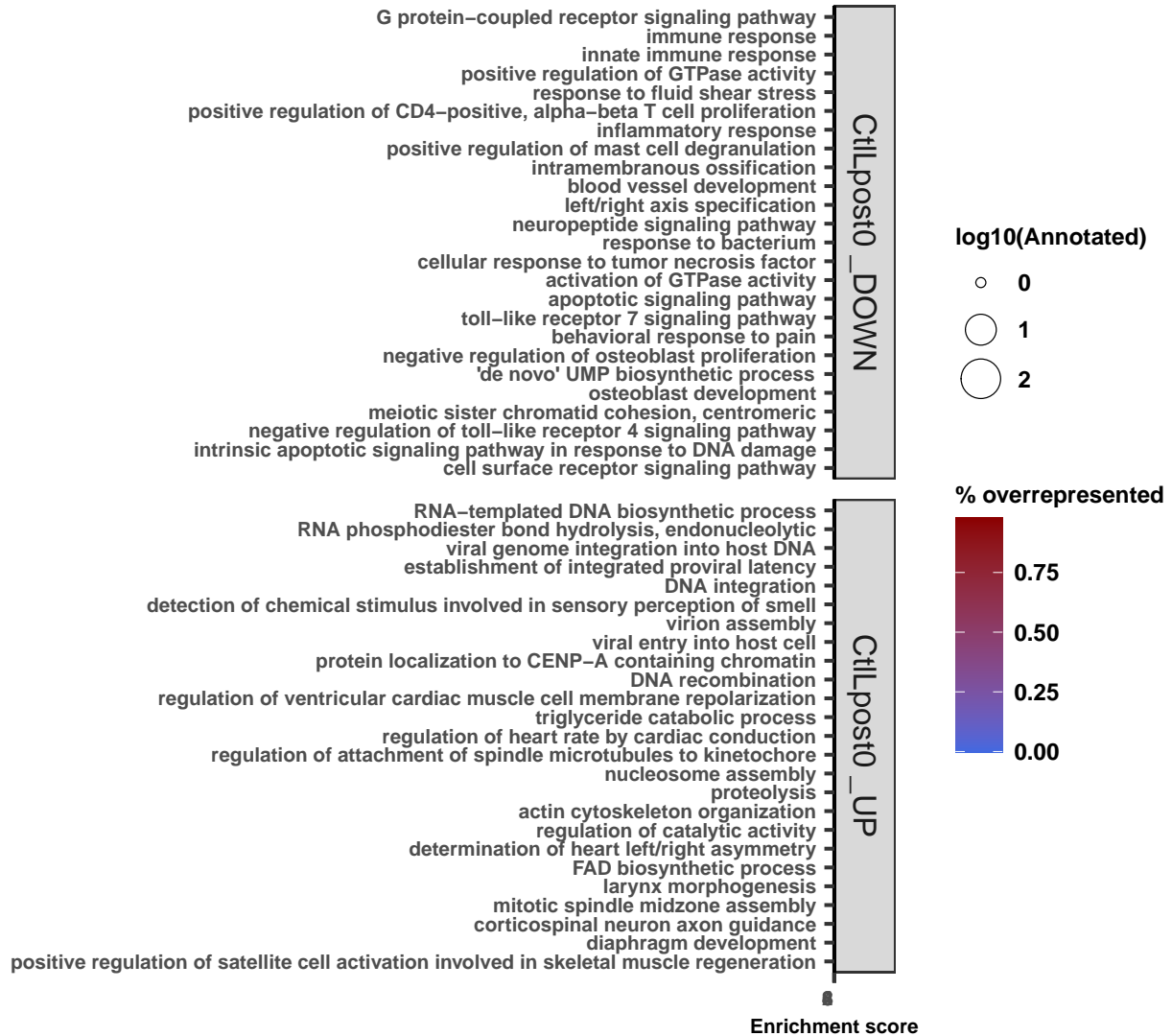
GO Biological processes upregulated in contrast: MGDietpost0

Top 25 terms ordered by weighted Fisher



GO Biological processes upregulated in contrast: CtlLpo:

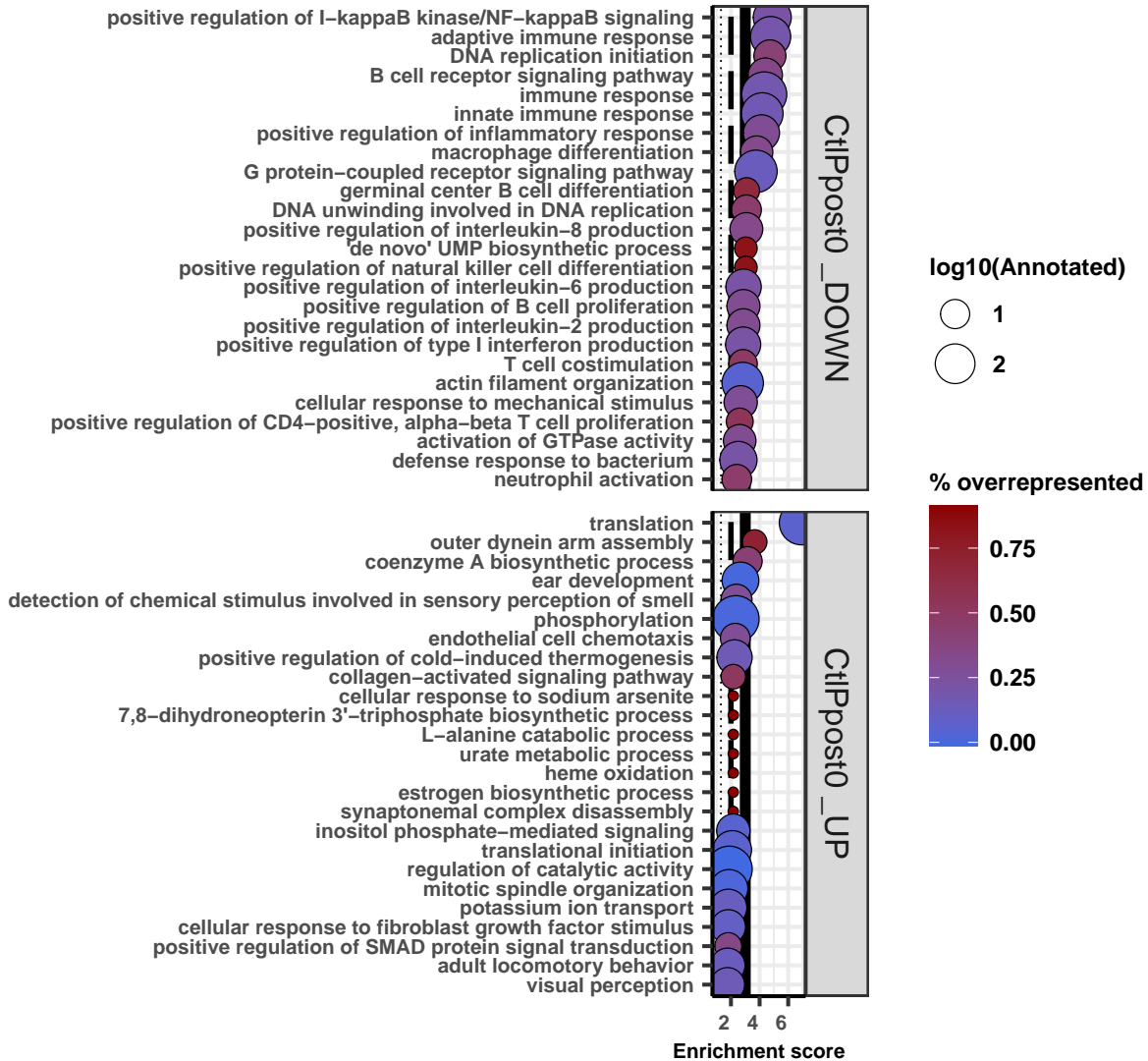
Top 25 terms ordered by weight



Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001

GO Biological processes upregulated in contrast: CtlPpost0

Top 25 terms ordered by weighted Fisher



```
# # Run GO enrichment on clustered genes from cluster_chunks.csv
#
# cluster_chunks_GO <- read_csv(here("cluster_analysis/cluster_chunks_GO.csv"))
#
# # Numeric ID values (in edgeR output) of genes whose names match Gene Ontology Estimates with edgeR
# chunk_matches <- which(toupper(names(canary_geneID2GO_BM)) %in% cluster_chunks_GO$gene)
# # these names came from `edgeR diff expression.R` file
#
# genesMatched <- unique(toupper(names(canary_geneID2GO_BM))[chunk_matches])
#
# DEgenes <- pull(cluster_chunks_GO, gene)
#
# # genes names with GO terms and DE.
# DEgeneswithGO <- which(DEgenes %in% genesMatched)
```

```

#
#   # select only genes that are DE and have GO terms
# DEgenesGO <- DEgenes[DEgeneswithGO]
#
#
# all_genes <- factor(as.integer(toupper(cluster_chunks_GO$gene) %in% DEgenesGO)) # create 0 1 factors
# names(all_genes) <- toupper(cluster_chunks_GO$gene) # add gene names
#
# #Do the GO mapping
# GOdata=new('topGOdata',
#             ontology='BP', # 'BP' = Biological Process. Also tried 'ALL' but doesn't work.
#             allGenes = all_genes, # all 10188 genes that are in MGpost0 data from edgeR. factor=1 for
#             annot = annFUN.gene2GO, # annotation function for custom annotation
#             # gene2GO = canary_geneID2GO_BM # annotation db from eggNOG output.
#             # I am thinking about switching this for one filtered to just edgeR genes:
#             gene2GO = canary_geneID2GO_BM_edgeR # annotation db from BiomaRt.
#             # gene2GO = canary_geneID2GO_BM_DE#,
#             # geneSel = DE_genes
#             )
#
# # Now analyze the results:
#
# # define test using the weight01 algorithm (default) with fisher
# weight_fisher_result=runTest(GOdata, algorithm='weight01', statistic='fisher')
#
#
# # generate a table of results: we can use the GenTable function to generate
# # a summary table with the results from tests applied to the topGOdata object.
# allGO=usedGO(GOdata)
# all_res=GenTable(GOdata,
#                  weightFisher=weight_fisher_result,
#                  orderBy='weightFisher',
#                  topNodes=length(allGO),
#                  numChar = 200)
#
# #performing BH correction on our p values
# p.adj=round(p.adjust(all_res$weightFisher,method="BH"),digits = 4)
#
# # create the file with all the statistics from GO analysis
# all_res_final=cbind(all_res,p.adj)
# all_res_final=all_res_final[order(all_res_final$p.adj),] %>%
#   mutate(across(everything(), gsub, pattern = ",", replacement = ";"))
#
# #get list of significant GO before multiple testing correction
# results.table.p= all_res_final[which(as.numeric(all_res_final$weightFisher)<=0.1),] %>%
#   mutate(across(everything(), gsub, pattern = ",", replacement = ";"))
#
#
# #get list of significant GO after multiple testing correction
# results.table.bh=all_res_final[which(all_res_final$p.adj<=0.1),] %>%
#   mutate(across(everything(), gsub, pattern = ",", replacement = ";"))
#
#

```

```
# #save first top 50 ontologies sorted by adjusted pvalues
# write.table(all_res_final[1:50,],
#             file = paste(
#               "~/Documents/GitHub/CanarySeq/geneAnnotations/geneID2GO_",
#               "cluster_chunks", ".csv", sep= ""),
#             # "~/Desktop/canary/summary_topGO_analysis.csv",
#             sep=";",
#             quote=FALSE,
#             row.names=FALSE)
#
#
```

```
library(scales)
```

```
ntop <- 25
```

```
ggdata_up <- all_res_up[1:ntop,] %>%
  mutate(enrichment = "lipid")
ggdata_up$Term <- factor(ggdata_up$Term, levels = rev(ggdata_up$Term)) # fixes order
```

```
ggdata_down <- all_res_down[1:ntop,] %>%
  mutate(enrichment = "protein")
ggdata_down$Term <- factor(ggdata_down$Term, levels = rev(ggdata_down$Term)) # fixes order
```

```
#
# gg1 <- ggplot(ggdata_up,
# # aes(x = Term, y = -log10(as.numeric(weightFisher)), size = -log10(as.numeric(weightFisher)), fill
# # aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = ((Significant-E
#
# geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
#           linetype = c("dotted", "longdash", "solid"),
#           colour = c("black", "black", "black"),
#           size = c(0.5, 1.5, 3)) +
#
# expand_limits(y = 1) +
# geom_point(shape = 21) +
# scale_size(range = c(2.5,12.5)) +
# scale_fill_continuous(low = 'royalblue', high = 'red4') +
#
# xlab('') +
# ylab('Enrichment score') +
# labs(
#   title = 'GO Biological processes upregulated in infected + lipid',
#   subtitle = 'Top 25 terms ordered by weighted Fisher p-value',
#   caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +
#
# theme_bw(base_size = 24) +
# theme(
#   legend.position = 'right',
#   legend.background = element_rect(),
#   plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust=0.5),
#   plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
#   plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),
```

```

#
#   axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
#   axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
#   axis.title = element_text(size = 12, face = 'bold'),
#   axis.title.x = element_text(size = 12, face = 'bold'),
#   axis.title.y = element_text(size = 12, face = 'bold'),
#   axis.line = element_line(colour = 'black'),
#
#   #Legend
#   legend.key = element_blank(), # removes the border
#   legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend
#   legend.text = element_text(size = 14, face = "bold"), # Text size
#   title = element_text(size = 14, face = "bold")) +
#
#   coord_flip()
#
#
#
# ggdata_down <- all_res_down[1:ntop,] %>%
#   mutate(enrichment = "protein")
# ggdata_down$Term <- factor(ggdata_down$Term, levels = rev(ggdata_down$Term)) # fixes order
#
# gg2 <- ggplot(ggdata_down,
#   aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = ((Significant-E
#
#   expand_limits(y = 1) +
#   geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
#     linetype = c("dotted", "longdash", "solid"),
#     colour = c("black", "black", "black"),
#     size = c(0.5, 1.5, 3)) +
#   geom_point(shape = 21) +
#   scale_size(range = c(2.5,12.5)) +
#   scale_fill_continuous(low = 'royalblue', high = 'red4') +
#
#   xlab('') +
#   ylab('Enrichment score') +
#   labs(
#     title = 'GO Biological processes upregulated in infected + protein',
#     subtitle = 'Top 25 terms ordered by weighted Fisher p-value',
#     caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +
#
#   theme_bw(base_size = 24) +
#   theme(
#     legend.position = 'right',
#     legend.background = element_rect(),
#     plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust = 0.5),
#     plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
#     plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),
#
#     axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
#     axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
#     axis.title = element_text(size = 12, face = 'bold'),
#     axis.title.x = element_text(size = 12, face = 'bold'),

```

```

#   axis.title.y = element_text(size = 12, face = 'bold'),
#   axis.line = element_line(colour = 'black'),
#
#   #Legend
#   legend.key = element_blank(), # removes the border
#   legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend
#   legend.text = element_text(size = 14, face = "bold"), # Text size
#   title = element_text(size = 14, face = "bold")) +
#
#   coord_flip()

# ggarrange(gg1, gg2, ncol=1)

ggdata <- bind_rows(ggdata_up, ggdata_down) %>%
  mutate(`% overrepresented` = (Significant-Expected)/Annotated)

gg <- ggplot(ggdata,
  aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = `% overrepresented`)) +
  expand_limits(y = 1) +
  geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
    linetype = c("dotted", "longdash", "solid", "dotted", "longdash", "solid"),
    colour = "black",
    size = c(0.5, 1.5, 3, 0.5, 1.5, 3)) +
  geom_point(shape = 21) +
  scale_size(range = c(2.5,12.5)) +
  scale_fill_continuous(low = 'royalblue', high = 'red4') +

  xlab('') +
  ylab('Enrichment score') +
  labs(
    title = 'G0 Biological processes upregulated by diet during infection',
    subtitle = 'Top 25 terms ordered by weighted Fisher p-value',
    caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +

  theme_bw(base_size = 24) +
  theme(
    legend.position = 'right',
    legend.background = element_rect(),
    plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust = 0.5),
    plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
    plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),

    axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
    axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
    axis.title = element_text(size = 12, face = 'bold'),
    axis.title.x = element_text(size = 12, face = 'bold'),
    axis.title.y = element_text(size = 12, face = 'bold'),
    axis.line = element_line(colour = 'black'),

    #Legend
    legend.key = element_blank(), # removes the border
    legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend

```

gg

Top 25 terms order



22

```

ggdata.c <- bind_rows(ggdata_up.c, ggdata_down.c) %>%
  mutate(`% overrepresented` = (Significant-Expected)/Annotated) %>%
  group_by(enrichment)

ggdata.c$Term <- factor(ggdata.c$Term, levels = unique(rev(ggdata.c$Term)), ordered=TRUE) # fixes order

gg.c <- ggplot(ggdata.c,
  aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = `% overrepresented`))

  expand_limits(y = 1) +
  geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
    linetype = c("dotted", "longdash", "solid", "dotted", "longdash", "solid"),
    colour = "black",
    size = c(0.5, 1.5, 3, 0.5, 1.5, 3)) +
  geom_point(shape = 21) +
  scale_size(range = c(2.5, 12.5)) +
  scale_fill_continuous(low = 'royalblue', high = 'red4') +

  xlab('') +
  ylab('Enrichment score') +
  labs(
    title = 'G0 Biological processes upregulated in Day 0 samples by Diet',
    subtitle = 'Top 25 terms ordered by weighted Fisher p-value',
    caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +

  theme_bw(base_size = 24) +
  theme(
    legend.position = 'right',
    legend.background = element_rect(),
    plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust = 0.5),
    plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
    plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),

    axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
    axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
    axis.title = element_text(size = 12, face = 'bold'),
    axis.title.x = element_text(size = 12, face = 'bold'),
    axis.title.y = element_text(size = 12, face = 'bold'),
    axis.line = element_line(colour = 'black'),

    #Legend
    legend.key = element_blank(), # removes the border
    legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend
    legend.text = element_text(size = 14, face = "bold"), # Text size
    title = element_text(size = 14, face = "bold")) +

  facet_grid(rows = vars(enrichment), scales = "free") +
  coord_flip()

gg.c

```


Day 1 Q/A



2

% overr

0.7

0.5

0.2

0.0

S. Look


```

ggdata.cpost0$Term <- factor(ggdata.cpost0$Term, levels = unique(rev(ggdata.cpost0$Term)), ordered=TRUE)

gg.cpost0 <- ggplot(ggdata.cpost0,
  aes(x = Term, y = -log10(as.numeric(weightFisher)), size = log10(Annotated), fill = `% overrepresented`)) +

  expand_limits(y = 1) +
  geom_hline(yintercept = c(-log10(0.05), -log10(0.01), -log10(0.001)),
    linetype = c("dotted", "longdash", "solid", "dotted", "longdash", "solid"),
    colour = "black",
    size = c(0.5, 1.5, 3, 0.5, 1.5, 3)) +
  geom_point(shape = 21) +
  scale_size(range = c(2.5, 12.5)) +
  scale_fill_continuous(low = 'royalblue', high = 'red4') +

  xlab('') +
  ylab('Enrichment score') +
  labs(
    title = 'GO Biological processes upregulated in Day 0 samples by Diet',
    subtitle = 'Top 25 terms ordered by weighted Fisher p-value',
    caption = 'Cut-off lines drawn at equivalents of p=0.05, p=0.01, p=0.001') +

  theme_bw(base_size = 24) +
  theme(
    legend.position = 'right',
    legend.background = element_rect(),
    plot.title = element_text(angle = 0, size = 16, face = 'bold', vjust = 1, hjust = 0.5),
    plot.subtitle = element_text(angle = 0, size = 14, face = 'bold', vjust = 1),
    plot.caption = element_text(angle = 0, size = 12, face = 'bold', vjust = 1),

    axis.text.x = element_text(angle = 0, size = 12, face = 'bold', hjust = 1.10),
    axis.text.y = element_text(angle = 0, size = 12, face = 'bold', vjust = 0.5),
    axis.title = element_text(size = 12, face = 'bold'),
    axis.title.x = element_text(size = 12, face = 'bold'),
    axis.title.y = element_text(size = 12, face = 'bold'),
    axis.line = element_line(colour = 'black'),

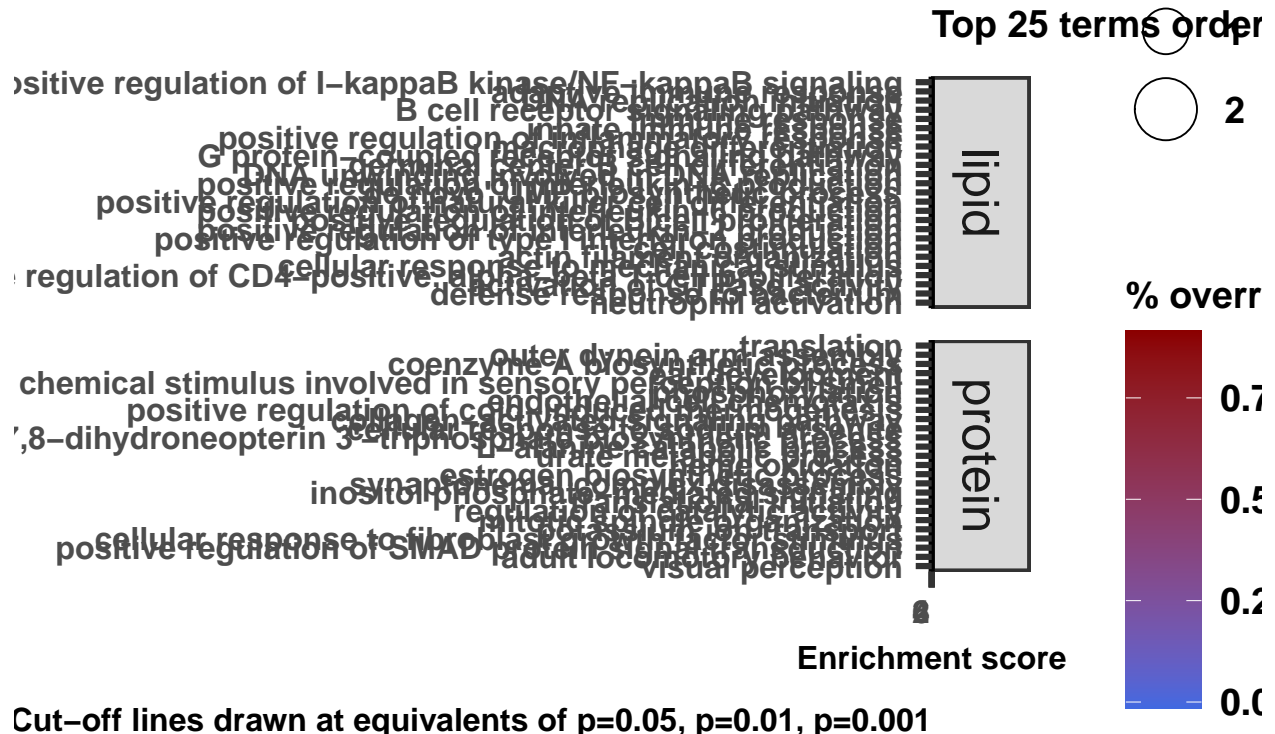
    #Legend
    legend.key = element_blank(), # removes the border
    legend.key.size = unit(1, "cm"), # Sets overall area/size of the legend
    legend.text = element_text(size = 14, face = "bold"), # Text size
    title = element_text(size = 14, face = "bold")) +

  facet_grid(rows = vars(enrichment), scales = "free") +
  coord_flip()

gg.cpost0

```

GO Biological processes upregulated in Day 0 As



Get genes all genes with a given GO term

```
# modify for easier formatting
ID2GO_up <- GOdata_up %>%
  genesInTerm() %>%
  enframe()

# find out how many contain GOTERM.
GOTERM <- "GO:0007186" #G Protein Coupled Receptor signaling pathway
GOTERM <- "GO:0002250" # adaptive immune response
GOTERM_ <- "GO:0030890" # positive regulation of B cell proliferation

# # select rows with GOTERM
ID2GO_up %>%
  filter(str_detect(name, GOTERM))

## # A tibble: 1 x 2
##   name      value
##   <chr>      <list>
## 1 GO:0002250 <chr [116]>
```

```

GcoupledUP <- genesInTerm(GOdata_up)[GOTERM][[1]][genesInTerm(GOdata_up)[GOTERM][[1]] %in% DEgenes_up]
GcoupledDN <- genesInTerm(GOdata_down)[GOTERM][[1]][genesInTerm(GOdata_down)[GOTERM][[1]] %in% DEgenes_down]

AdImUP <- genesInTerm(GOdata_up)[GOTERM][[1]][genesInTerm(GOdata_up)[GOTERM][[1]] %in% DEgenes_up]
AdImDN <- genesInTerm(GOdata_down)[GOTERM][[1]][genesInTerm(GOdata_down)[GOTERM][[1]] %in% DEgenes_down]

# b cell pos reg (Diet post0)
PosBUP <- genesInTerm(GOdata_up)[GOTERM_][[1]][genesInTerm(GOdata_up)[GOTERM_][[1]] %in% DEgenes_up]

# let's run the same thing but not for a given go term
# reformat genesInTerm as vector
all_up_genesGOTERM <- names(unlist(genesInTerm(GOdata_up)))
all_down_genesGOTERM <- as.vector(unlist(genesInTerm(GOdata_down)))

allUP <- all_up_genesGOTERM[all_up_genesGOTERM %in% DEgenes_up]
allDN <- all_down_genesGOTERM[all_down_genesGOTERM %in% DEgenes_down]

Gcoupled <-
  list(
    PROTEIN = GcoupledUP,
    LIPID = GcoupledDN
  )

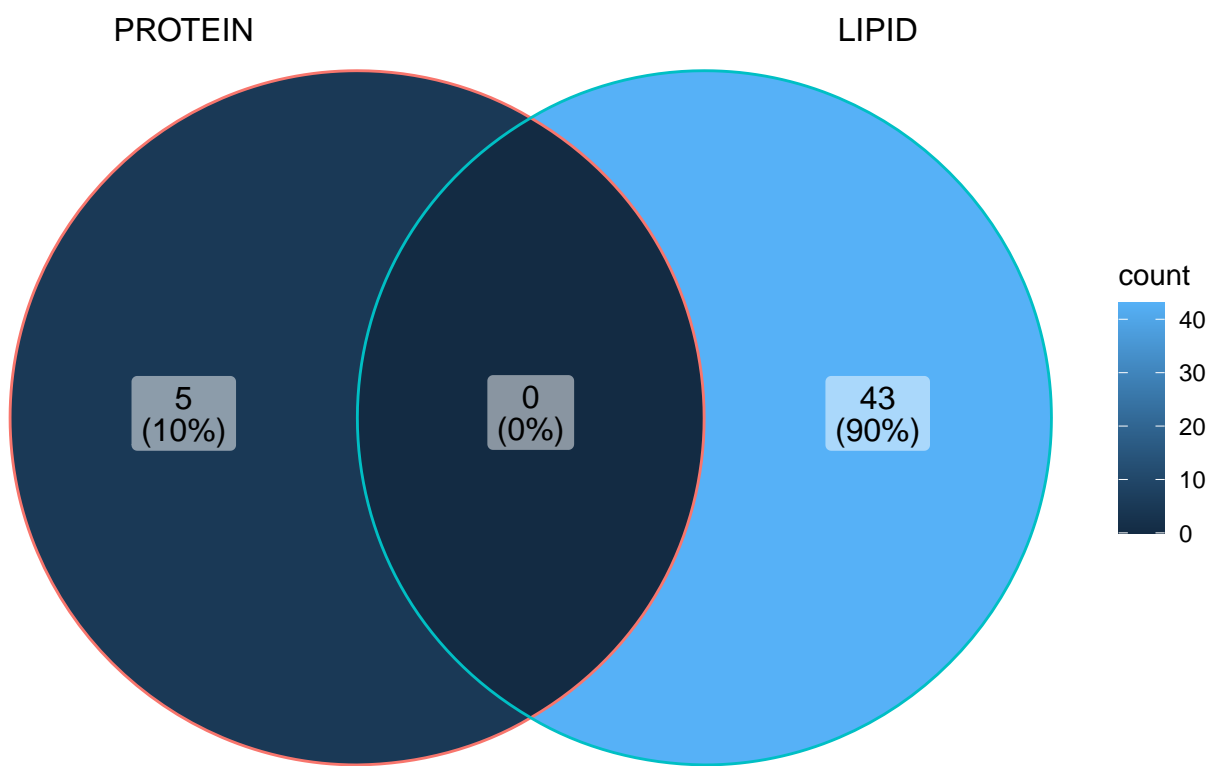
AdaptImm <-
  list(
    PROTEIN = AdImUP,
    LIPID = AdImDN
  )

all <-
  list(
    PROTEIN = allUP,
    LIPID = allDN
  )

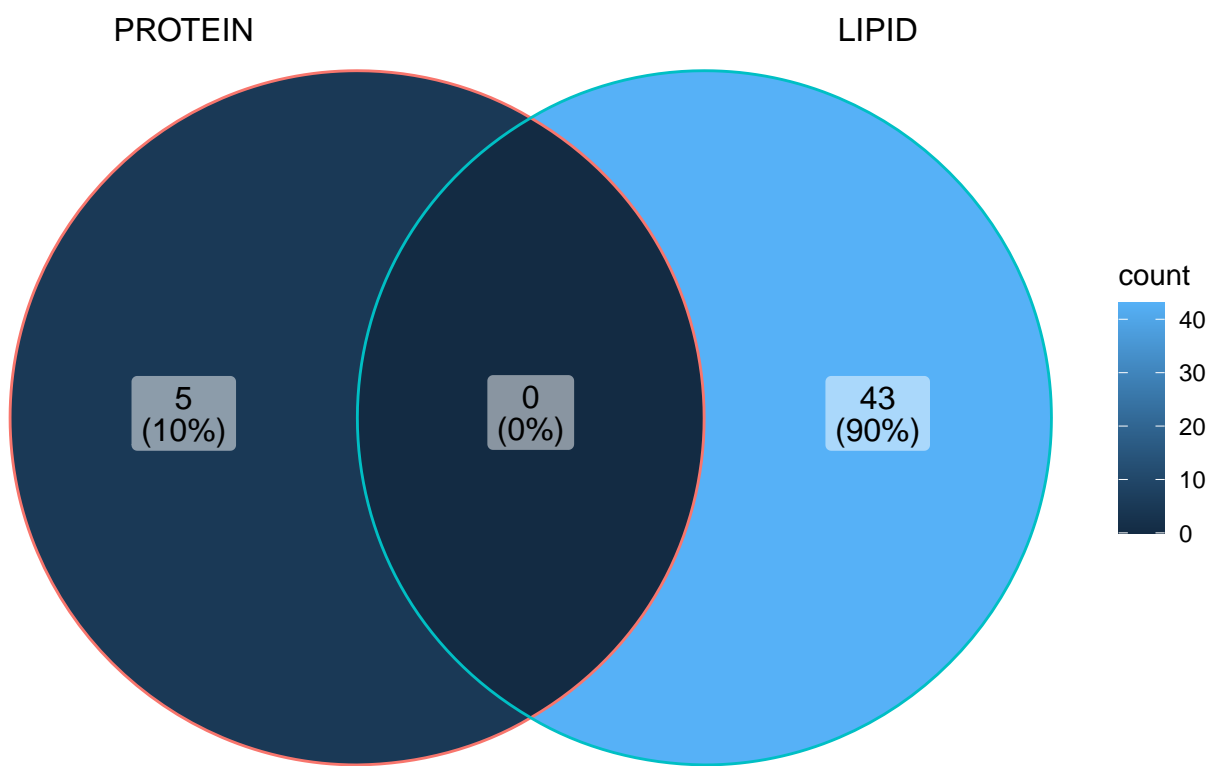
# if (!require(devtools)) install.packages("devtools")
# devtools::install_github("gaospecial/ggVennDiagram")
library(ggVennDiagram)

ggVennDiagram(Gcoupled)

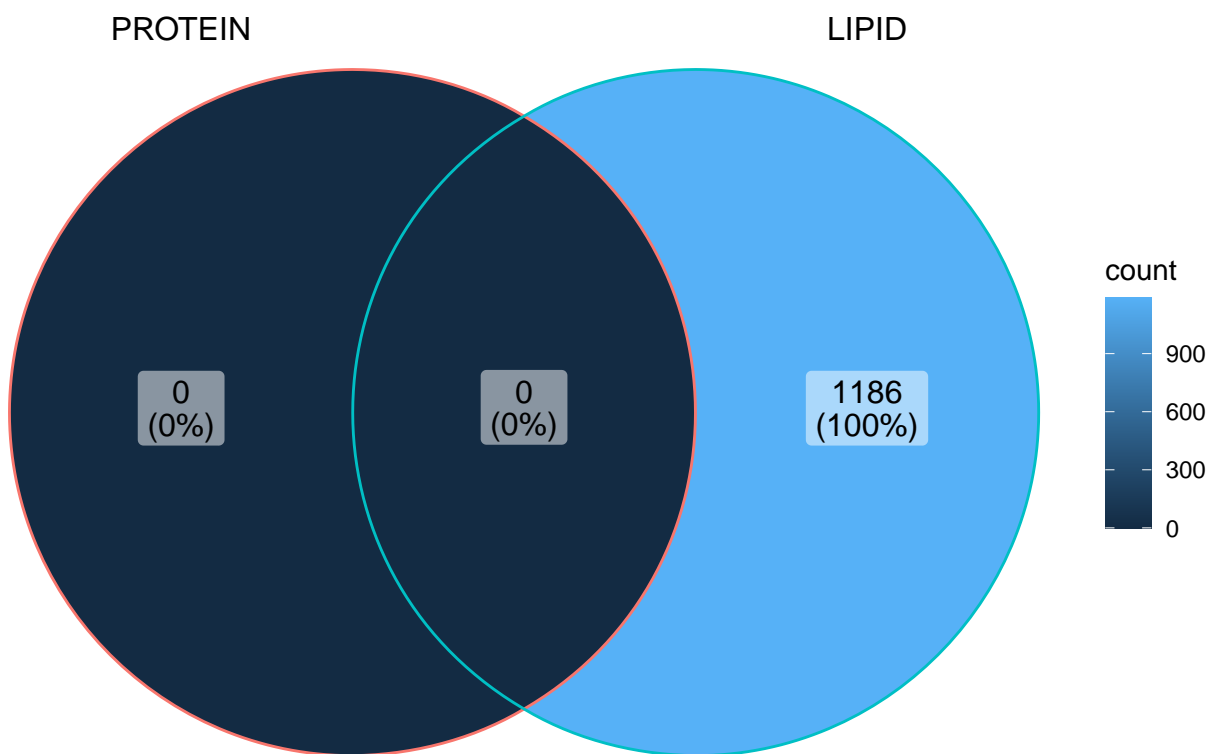
```



```
ggVennDiagram(AdaptImm)
```



```
ggVennDiagram(all)
```



```
# ADA
# AGTRAP
# AKAP13
# C3
# CCR10
# CCR6
# CCR7
# DGKH
# DGKZ
# GNG2
# GNG7
# GPR132
# GPR162
# GPR171
# GPR174
# GPR18
# GPR82
# LOC103812291
# LOC103812980
# LOC103812985
# LOC103813169
# LOC103813540
# LOC103815708
# LOC103819762
# LOC103824208
# LOC108964377
```

LOC108964492
LOC108964818
LOC115485368
LOC115485604
LOC115485605
P2RY10
P2RY12
P2RY6
PIK3CG
PIK3R6
PLCL1
PLCL2
PREX1
PRMT5
PTK2B
RAMP2
RGS1
RGS4
RRH
S1PR2
S1PR4
SLC24A4
SORL1
TYRO3
ZDHHC21