

Additional File 2: Scripts used in this study

```
1
2
3 #Pearson's coefficient correlation with complete hierarchical agglomerative clustering
4
5 library(factoextra)
6
7
8 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
9
10 inData <- "FILENAME"
11 df <- data.frame(read.csv(paste0(inData, ".csv"),
12                             header=T,
13                             row.names=1,
14                             check.names=FALSE))
15
16 df_scaled <- scale(df, center = T, scale = T)
17
18
19 dis.cor <- get_dist(df_scaled, method = 'pearson')
20
21 clusters <- hclust(dis.cor, method = 'complete')
22
23 i=2
24 while (i < 41){
25     clust <- cutree(clusters, k=i)
26     clust <- as.data.frame(clust)
27     df$cluster <- clust$clust
28
29     write.csv(df, paste0('Pearsons_', i, '_clusters.csv'))
30     print(i)
31     i = i+1
32 }
33
34 #weighting the variables
35 t_df_scaled <- t(df_scaled)
36
37 weighted_corr <- cov.wt(t_df_scaled, wt = c(6,6,1,1,1,1), cor = T) 38
39 corr.matrix <- weighted_corr$cor
40
41 corr.dis <- dist(corr.matrix)
42
43 corr.data <- hclust(corr.dis, method = 'complete')
44
45 i=2
46 while (i < 41){
47     clust <- cutree(corr.data, k=i)
48     clust <- as.data.frame(clust)
49     df$cluster <- clust$clust
50
51     write.csv(df, paste0('Weighted_pearsons_', i, '_clusters.csv'))
52     print(i)
53     i = i+1
54 }
55
```

```

1 #K-means clustering
2
3 library(FactoMineR)
4 library(factoextra)
5 library(cluster)
6
7 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
8 inData <- 'FILENAME'
9
10 data <- read.csv(paste0(inData, '.csv'),
11                 header=T,
12                 sep=";",
13                 row.names=1,
14                 check.names=FALSE)
15 data <- as.matrix(data)
16 i=2
17
18 while (i< 50) {
19   kdata <- kmeans(data, center = i, nstart=25)
20
21   gap_stat <- clusGap(data, FUN=kmeans, nstart =25, K.max=i, B=50)
22   print(i)
23
24   distance <- get_dist(data)
25
26   #creates a pdf of several plots
27   pdfname <- paste0('K-means',i,'clusters.pdf')
28   pdf(pdfname,
29       width = 11,
30       height = 11)
31
32   fviz_nbclust(data, kmeans, method='wss', k.max=i)
33   fviz_nbclust(data, kmeans, method='silhouette', k.max=i)
34   fviz_gap_stat(gap_stat)
35   fviz_dist(distance)
36   fviz_cluster(kdata, data=data)
37
38   dev.off()
39
40   #creates spreadsheet of input data with cluster assignments
41   attach_Data$cluster <- kdata$cluster
42   csvname <- paste0('K-means',i,'clusters.csv')
43   write.csv(attach_Data, csvname)
44   i = i+1
45 }

```

```
1 #Hierarchical clustering on principle components
2
3 library(FactoMineR)
4 library(factoextra)
5
6 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
7
8 inData <- 'FILENAME'
9
10 data <- read.csv(file=paste0(inData, '.csv'),
11                 header=T,
12                 sep=',',
13                 row.names=1,
14                 check.names=FALSE)
15 data <- as.matrix(data)
16
17 just_data <- data[,6:15]
18 i=2
19 while (i< 41){
20   pca <- FAMD(data)
21   res.hcpc <- HCPC(pca,-1,iter.max = 50, min= i, max =i, graph =T)
22
23   #creates spreadsheet of input data with cluster assignments
24   csvname <- paste0('HCPC clustering ',i,'clusters.csv')
25   write.csv(res.hcpc$data.clust,csvname)
26   print(i)
27   i=i+1
28 }
29
```

```

1 #Partitioning around medoids
2 library(cluster)
3
4 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
5 #load data
6 inData <-
7 dat =read.csv("KO_Phenotypes_ALL_10_7_28_19_FINAL.csv")
8 #subset out measurements
9 dat1 = dat[,13:18]
10
11 #load library to creat dissimilarity matrix with mixed data types (gowers)
12 library(cluster)
13
14 #Different weights for variables
15 weight0 <- c(1,1,1,1,1,1,1,1,1,1)
16 weight1 <- c(2,2,1,1,1,1,1,1,1,1)
17 weight2 <- c(2,2,1,.5,1,.5,1,.5,1,.5)
18 weight3 <- c(6,2,2,2,1,1,1,1,1,1)
19 weight4 <- c(6,5,.5,.5,1,1,1,1,1,1)
20 weight5 <- c(6,4,1,1,1,1,1,1,1,1)
21 weight6 <- c(6,6,1,1,1,1,1,1,1,1)
22 #calculate gower's distance
23 distF = daisy (dat1, metric = 'gower', weight = weight0)
24
25 #Partitioning around medoids
26 i <- 2
27 while (i<51) {
28 pamfit = pam(distF, diss=TRUE, k=i)
29 #obtain clusters
30 clust = pamfit$clustering
31 ## add clusters to data
32 dat$cluster = clust
33 #export clusters
34 write.csv(dat, file = paste0('PAM_',i,'_clusters.csv'))
35 i = i+1
36 }
37
38 #Ward's minimum variance
39 hc_fun = hclust(distF, method="ward.D")
40
41 i=2
42 while (i<51){
43 clusters <- cutree(hc_fun, k=i)
44 data$cluster <- clusters
45 write.csv(data, paste0('Ward_',i,'_clusters.csv'))
46 print(i)
47 i=i+1
48 }
49
50
51

```

```
1 #gets the number per cluster
2 library(plyr)
3
4 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
5
6 i=2
7 statlist <- list()
8
9 while (i<41) {
10   name <- paste0('FILENAME',i, '.csv')
11   data1 <- read.csv(name)
12   data2 <- data1$ANY_COLUMN
13   Number <- as.data.frame(data2)
14   Number$cluster <- sheet$clust
15
16   number <- ddply(Number, ('cluster'), summarise,
17     N = length(data2))
18
19   statlist[[i-1]] <- number$Number
20
21   print(i)
22   i=i+1
23 }
24
25 COUNTED <- as.data.frame(t(ldply(statlist, rbind)))
26 write.csv(COUNTED, 'NEW CSV FILENAME.csv')
27
```

```

1 #Binning genes into categories per cluster
2 library(plyr)
3 library(dplyr)
4 library(cowplot)
5
6 filename <- 'FILENAME.csv'
7 data <- read.csv(data)
8
9 #Group the data by a factor (column)
10 data$group1 <- cut('Grouping factor1', c(0,40,65,75,77.5,82.5,85,100)) #Grouping factor
can be any column from data, numeric vector (c(...)) is provided as an example
11 data$group2 <- cut('Grouping factor2', c(0,15,25,30,35,40,45,100))
12
13 #Converts the numeric cuts into categories
14 levels(data$group1) <- c('Severly Reduced','Reduced','Slightly Reduced','Normal Low',
'Normal Average', 'Normal High', 'Increased')
15 levels(data$group2) <- c('Severly Reduced','Reduced','Slightly Reduced','Normal Low',
'Normal Average', 'Normal High', 'Increased')
16
17 supersheet <- data[, (ncol(data)-2):(ncol(data))]
18
19 sum <- supersheet %>%
20   group_by(cluster, group) %>%
21   summarise(number = length(group)) %>%
22   mutate(percent = (number / sum(number))*100)
23
24 sumah <- supersheet %>%
25   group_by(cluster, ahgroup) %>%
26   summarise(ahnumber = length(ahgroup)) %>%
27   mutate(percent = (ahnumber / sum(ahnumber))*100)
28
29 #Plots
30 Growth <- ggplot(sum, aes(x=cluster, y = percent, fill = group))+
31   geom_bar(stat = 'identity')+
32   geom_text(data = subset(sum,percent !=0),aes(label=paste0(round(percent), "%")),
33     position = position_stack(vjust = 0.5), size =4, fontface='bold')+
34   geom_text(data = sum, aes(label=number), position = position_stack(vjust=0.4), size =
35     3, fontface='bold')+
36   theme(axis.text.x = element_text(angle = 65, hjust = 1, face='bold', color='black',
37     size=20),
38     axis.text.y = element_text(face='bold', color='black', size=20),
39     axis.title = element_text(face='bold', color='black', size=20 ),
40     plot.title = element_text(face='bold', color='black', size=20))
41
42 AH <- ggplot(sumah, aes(x=cluster, y = percent, fill = ahgroup))+
43   geom_bar(stat = 'identity')+
44   geom_text(data = subset(sumah,percent !=0),aes(label=paste0(round(percent), "%")),
45     position = position_stack(vjust = 0.5), size =4, fontface='bold')+
46   geom_text(data = sumah, aes(label=ahnumber), position = position_stack(vjust=0.4),
47     size = 3, fontface='bold')+
48   theme(axis.text.x = element_text(angle = 65, hjust = 1, face='bold', color='black',
49     size=20),
50     axis.text.y = element_text(face='bold', color='black', size=20),
51     axis.title = element_text(face='bold', color='black', size=20 ),
52     plot.title = element_text(face='bold', color='black', size=20))
53
54 summary <- merge(sum, sumah, by='cluster')
55
56 write.csv(summary, 'PAM 0 20 Growth and AH summary.csv')
57
58 pdf('Binned groups by cluster PAM 0 20.pdf', height = 12, width = 36)
59 plot_grid(Growth, AH, labels =c('A.', 'B.'), nrow = 1)
60 dev.off()

```

```

1 #Average consensus and average relative standard deviation
2
3 library(plyr)
4 library(dplyr)
5 library(tidyr)
6 library(purrr)
7
8 #####FUNCTIONS#####
9 count_words <- function(data, column, clustercol) {
10   datacolumn <- data[,column]
11   datacolumn <- as.data.frame(datacolumn)
12   datacolumn$cluster <- data[,clustercol]
13
14   df <- datacolumn %>%
15     group_by(cluster)%>%
16     mutate(words= strsplit(as.character(datacolumn), ' '))%>%
17     unnest()%>%
18     dplyr::count(cluster, words) %>%
19     transmute(n, pcnt = (n/sum(n)*100))
20
21   df <- ddply(df, c('cluster'), summarise, value = max(pcnt))
22   max_pcnt_average <- mean(df$value)
23   return(max_pcnt_average)
24 }
25
26 relative_standard_deviations <- function(data, column, clustercol) {
27   df <- data[,column]
28   df <- as.data.frame(df)
29   df$cluster <- data[,clustercol]
30
31   summary_data <- ddply(df, c('cluster'), summarise,
32     N = length(df),
33     mean = mean(df),
34     sd = sd(df),
35     relative_sd = ((sd/mean)*100))
36   relative_sd_means <- colMeans(summary_data[,5,drop=FALSE])
37   return(relative_sd_means)
38 }
39 #####
40
41 setwd('YOUR_WORKING_DIRECTORY/FOLDER')
42
43 #Empty lists where information will be stored
44 rgsdlist <- list()
45 rahsdlist <- list()
46 rwordlist <- list()
47 rfinalwordlist <- list()
48
49 #The main loop that calculates the relative standard deviation of phenotypic traits per
50 cluster
51 i= 20
52 while (i<34) {
53   name <- paste0('FILENAME_',i,'_clusters.csv')
54
55   sheet <- read.csv(name, row.names = 1)
56   sheet <- sheet[,9:19] #choose columns that need to be analyzed
57
58   rsdmeans <- relative_standard_deviations(sheet, 1, 11)
59   sdname <- paste('Gsd',i)
60   rgsdlist[[sdname]] <- rsdmeans #adds the relative standard deviation for one
61   continuous variable to a list
62
63   rsdmeans2 <- relative_standard_deviations(sheet, 2, 11)
64   sdname <- paste('AHsd',i)
65   rahsdlist[[sdname]] <- rsdmeans2 #adds the relative standard deviation for one
66   continuous variable to a list
67 }

```

```

65     #iterates through the columns with categorical data
66     word_clusters <- paste0('cluster',i)
67     k=3
68     while (k < 11){
69         rwordlist[[k-2]] <- count_words(sheet, k, 11)
70         rfinalwordlist[[word_clusters]] <- rwordlist
71         k=k+1
72     }
73     print(i)
74     i=i+1
75 }
76
77 #Binds the data stored in lists to a dataframe and then outputs a .csv file
78 do.call('rbind', rfinalwordlist)
79 relative_sd_frame <- data.table::rbindlist(rfinalwordlist)
80 names(relative_sd_frame) <- c('C num', 'C morph', 'PP num', 'PP morph', 'P num', 'P
morph', 'A num', 'A morph') #names the columns
81 relative_sd_frame$growth <- unlist(rgsdlist)
82 relative_sd_frame$AH      <- unlist(rahsdlist)
83
84 write.csv(relative_sd_frame, 'Summary_info_PAM.csv')
85

```

```

1 #K-means clustering of expression data
2
3 my_scale <- function(x) { #scales the data between -2 and 2
4   m = apply(x, 1, mean, na.rm = T)
5   s = apply(x, 1, sd, na.rm = T)
6   return((x - m) / s)
7 }
8
9 wang <- factor(c(0,2,24,48,72,96,120,144))
10 wangnum <- c('0 h', '2 h', '24 h', '48 h', '72 h', '96 h', '120 h', '144 h')
11 germ <- factor(c(0.25,2,4,6))
12 germnum <- c('1/4 h', '2 h', '4 h', '6 h')
13 green <- factor(c(0,2,4,8,10,12,12.2,14,18,24))
14 greennum <- c('0 h', '2 h', '4 h', '8 h', '10 h', '12B', '12T', '14 h', '18 h', '24 h')
15 wanggerm <- factor(c(0.25, 2, 4, 6))
16 wanggermnum <- c('0.25 h', '2 h', '4 h', '6 h')
17 kasuga <- factor(c(1,3,9,15,21,27))
18 kasuganum <- c('0 hr', '3 hrs', '9 hrs', '15 hrs', '21 hrs', '27 hrs')
19
20 library(ggplot2)
21 library(ggthemes)
22 library(lemon)
23 library(cluster)
24 library(factoextra)
25
26 setwd('YOUR_WORKING_FOLDER')
27
28 inData <- 'FILENAME'
29 i <- 7 #number of clusters
30
31 df <- data.frame(read.csv(file=paste(inData, ".csv", sep=""),
32                       header=T,
33                       row.names=1,
34                       check.names=FALSE))
35 clean_df <- df[1:10] #columns with data for scaling
36 cleaner_df <- my_scale(clean_df)
37
38 fit <- kmeans(cleaner_df, i, nstart = 50)
39 output <- data.frame(cleaner_df, fit$cluster)
40 output <- as.data.frame(output)
41
42 graph <- tibble::rownames_to_column(output, 'gene')
43 cluster_graph <- tidyr::gather(graph, key, value, -fit.cluster, -gene)
44 uh <- cluster_graph$key
45 uh <- gsub("[a-zA-Z]", "", uh)
46 uh <- as.numeric(uh)
47 cluster_graph$key <- uh
48 cluster_graph$key <- as.factor(cluster_graph$key)
49 cluster_graph <- as.data.frame(cluster_graph)
50 cluster_graph$TF <- df$TF
51
52 tiff(paste0('k-means green clustering ', i, ' clusters.tiff'),
53      width=3840,
54      height=2160,
55      units='px')
56
57 ggplot(cluster_graph, aes(key, value, group = gene, ))+
58   geom_point(aes(color = factor(fit.cluster)))+
59   geom_line(aes(color = factor(fit.cluster)))+
60   facet_rep_wrap(~fit.cluster, repeat.tick.labels='all')+
61   theme_tufte()+
62   theme(axis.line = element_line(), legend.position = 'none',
63         axis.text.x = element_text(angle = 65, hjust = 1, face='bold',
64                                     color='black', size=40),
65         axis.text.y = element_text(face='bold', color='black', size=40),
66         axis.title = element_text(face='bold', color='black', size=40),
67         strip.text = element_text(face='bold', size=40))+
68   scale_color_manual(breaks = cluster_graph$fit.cluster, values = c('#e6194b',
69                             '#3cb44b', '#ffe119', '#4363d8', '#f58231', '#911eb4',

```

```
68         '#46f0f0', '#f032e6', '#bcf60c', '#fabebe', '#008080', '#e6beff',
69         '#9a6324', '#fffac8', '#800000', '#aaffc3', '#808000', '#ffd8b1',
70         '#000075', '#808080', '#ffffff', '#000000'))+
71     scale_x_discrete(limits = green,
72                     labels = greennum)
73
74
75 dev.off()
76
77 output$TF <- df$TF
78 write.csv(output, paste0('K-means green ', i, ' clusters with TF.csv'))
79
80
```