

R implementation

Sleiman Bassim, PhD

September 7, 2016

1 Loaded functions:

```
#source("/media/Data/Dropbox/humanR/01funcs.R")
rm(list=ls())
#setwd("/media/Data/Dropbox/humanR/PD/")
#setwd("~/Dropbox/humanR/PD/")
###load("PD.Rdata", .GlobalEnv)
#lsos(pat="")
```

2 Load packages.

```
pkgs <- c('xlsx', 'caret', 'leaps', 'glmnet', 'lattice',
          'latticeExtra', 'dplyr', 'tidyr', 'grid',
          'gplots', 'WGCNA', 'reshape2', 'stringr')
lapply(pkgs, require, character.only = TRUE)
```

3 1 Summary of the workflow

For pathway construction we will be using a full featured pipeline that implements i- gene network inference from RNA-seq gene expression data, ii- gene discovery after protein domain alignments to ten well cited protein databases, iii- gene interaction validation from inferred networks published in the European private String database, and finally iv- regularization analysis to eliminate to least probabilistic networks identified, however are less likely to be accurately constructed. Briefly, R script is used to calculate a weighted matrix for every gene discovered from the RNA-seq data of ganglia samples. This correlated matrix will be used to select genes with similar patterns of expression that will constitute an ensemble of gene modules. Each module includes eigengenes, which through the guilt-by-association method are identified to be signatures of a cell response at a precise moment after detection of a stimuli. These group of genes will go through a second set of analysis to calculate a probability of interaction at different error thresholds. Next step is to identify the actual genes for their protein functions, hence create an image of the interactions between biological processes specific to these functions. Moreover, validation of the gene networks will depend on the nature of these processes. There is many online sequence databases that deliver this information after acquiring gene names and protein functions. However, we will be using an approach that contains both statistical inference of gene interactions and functional association between processes. Therefore, after constructing gene-gene interaction networks we will compare them to already established networks inferred from all published data over many clades and between several closely related species. Finally, we will use machine learning methods, mainly regularization methods which are known to reduce the hyperspace of the data (that is reduce the number of candidates, eg., genes, networks) and let us focus on the most probabilistically accurate results by keeping the least biased gene networks. Overall the pipeline has already been tested on different datasets, databases are already been acquired and deployed on our servers, and the R code and many of the packages to be used have been tested and some published.

27 2 Differentially expressed genes

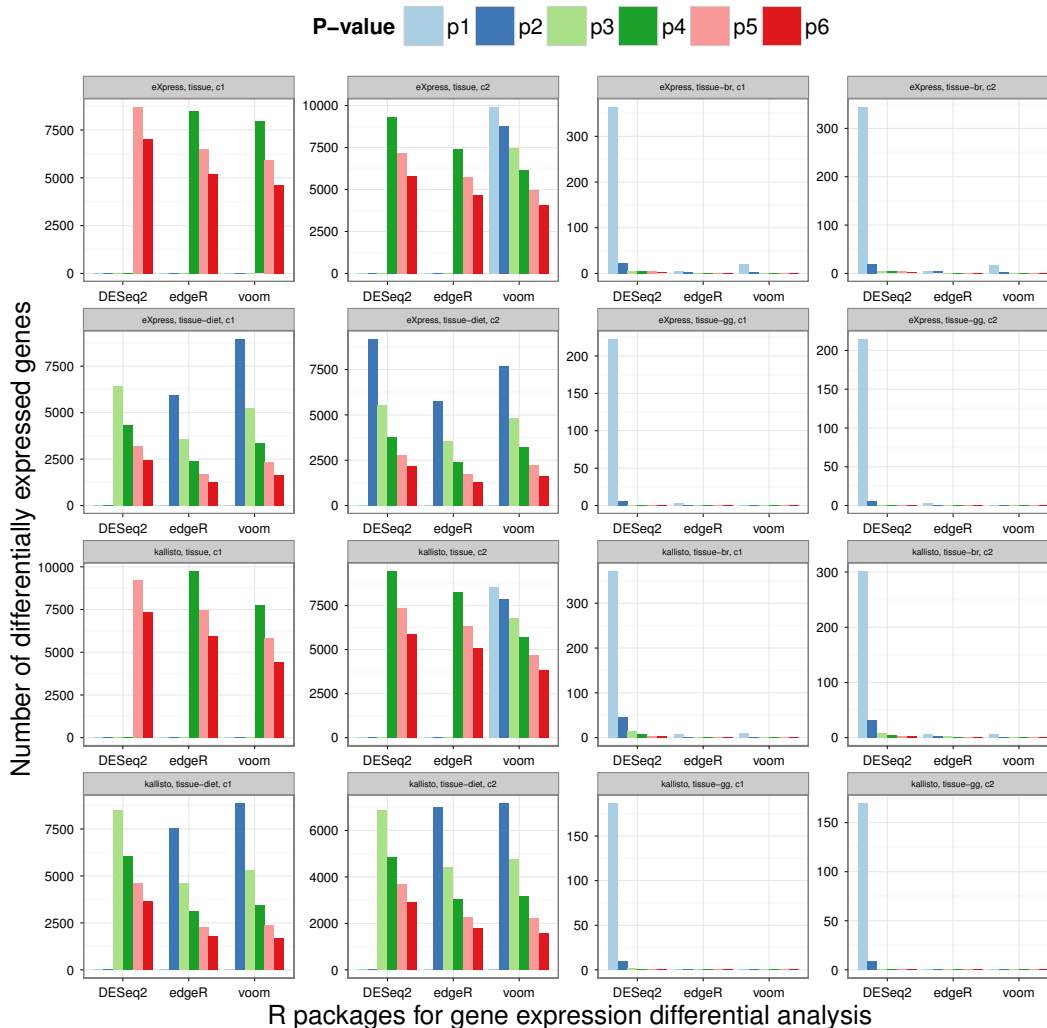
Differentially expressed genes are counted from mapping **both gills and ganglia** sequenced samples to reference transcriptome built from all samples.

```
read.table("./data/summary.raw.all.txt") %>%
```

```

ggplot(aes(
  x = V1,
  y = V8,
  fill = V6)) +
theme_bw() +
geom_bar(stat = "identity",
  position = "dodge") +
facet_wrap(~ V4 + V5 + V7,
  ncol = 4,
  scales = "free") +
scale_fill_brewer(type = "qual", palette = "Paired",
  name = "P-value") +
labs(x = "R packages for gene expression differential analysis",
  y = "Number of differentially expressed genes") +
theme(legend.position = "top",
  axis.text.x = element_text(vjust = .5, size = 6),
  axis.text.y = element_text(vjust = .5, size = 6),
  strip.text = element_text(size = 4),
  axis.ticks = element_line(size = .2),
  axis.ticks.length = unit(.05, "cm"))

```



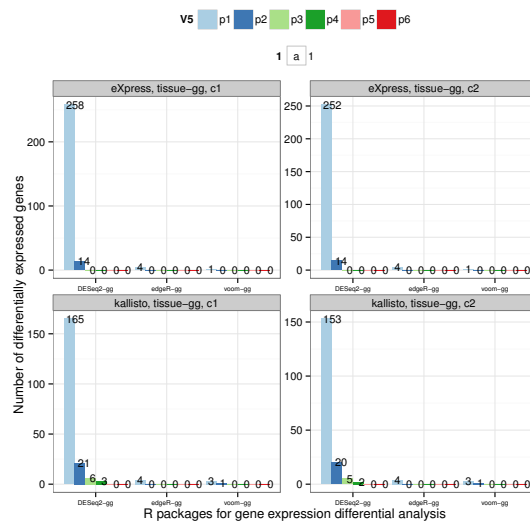
30
31 Differentially expressed genes are counted from mapping **ganglia** sequenced samples to reference tran-
32 scriptome built from all samples.

```
read.table("./data/summary.gg.txt") %>%
```

```

ggplot(aes(
  x = V2,
  y = V8,
  fill = V5)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V2,
    y = V8,
    ymax = V8,
    label = V8,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V3 + V4 + V6,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



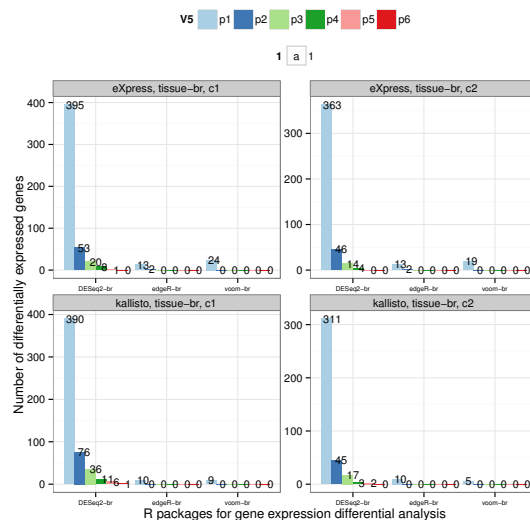
Differentially expressed genes are counted from mapping **gills** sequenced samples to reference transcriptome built from all samples.

```
read.table("./data/summary.br.txt") %>%
```

```

ggplot(aes(
  x = V2,
  y = V8,
  fill = V5)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V2,
    y = V8,
    ymax = V8,
    label = V8,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V3 + V4 + V6,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



2.1 Increasing DEG by changing the trimming rates of raw reads

Getting gene expression by mapping the original raw reads **without trimming** to the gills de novo transcriptome.

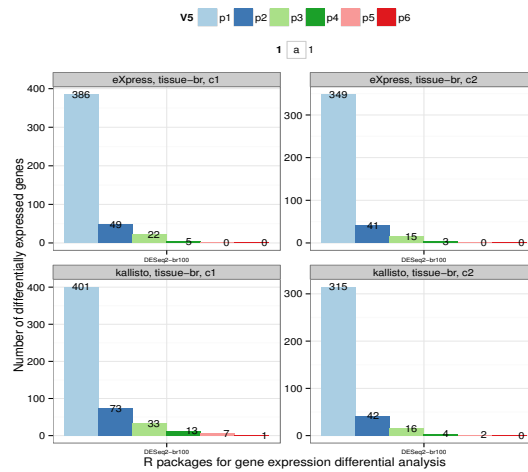
† De novo assembly was carried out with trimmed reads though

```
read.table("./data/summary.br_35078.txt") %>%
```

```

ggplot(aes(
  x = V2,
  y = V8,
  fill = V5)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V2,
    y = V8,
    ymax = V8,
    label = V8,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V3 + V4 + V6,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



2.2 Increasing DEGs by changing the normalization strategy: Fast abundance quantification *kallisto*

The below graph shows the number of differentially expressed genes when raw reads were normalized separately for each biological sample.

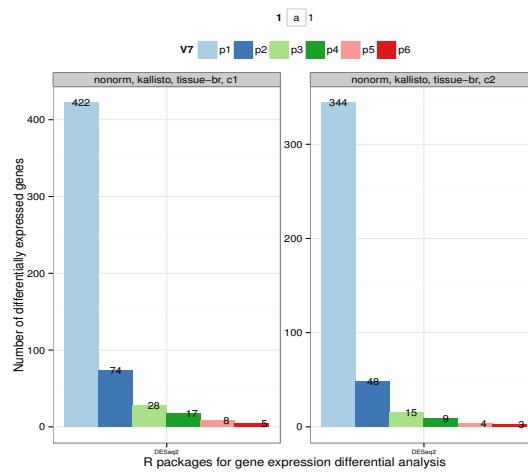
† All the analyses before were done on normalized reads by grouping all biological samples together

```
read.table("./data/summary.br.nonorm.txt") %>%
```

```

ggplot(aes(
  x = V1,
  y = V10,
  fill = V7)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V1,
    y = V10,
    ymax = V10,
    label = V10,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V4 + V5 + V6 + V8,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = 'Paired') +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



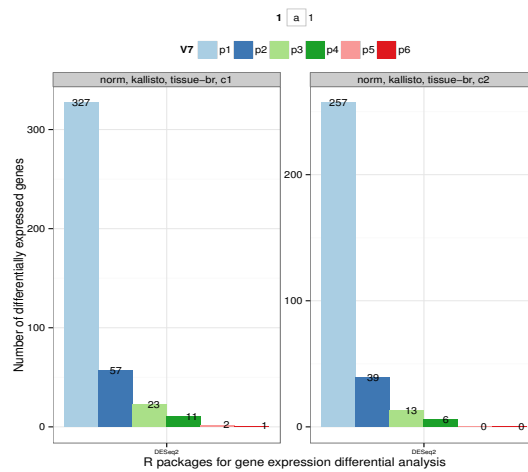
The below graph shows the number of differentially expressed genes when raw reads were **NOT** normalized.

```
read.table("./data/summary.br.norm.txt") %>%
```

```

ggplot(aes(
  x = V1,
  y = V10,
  fill = V7)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V1,
    y = V10,
    ymax = V10,
    label = V10,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V4 + V5 + V6 + V8,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



2.3 Increasing DEGs by changing the normalization strategy: abundance quantification with alignment *express*

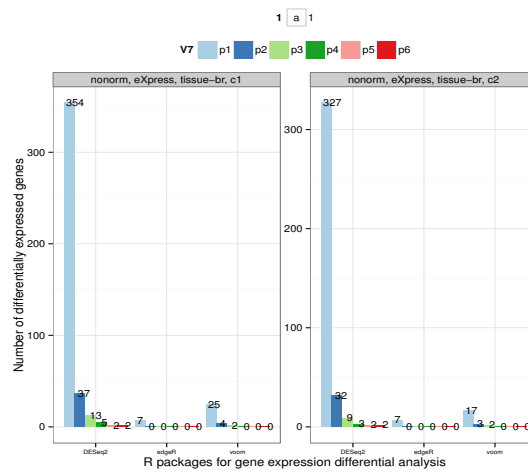
The R package *eXpress* is solely used to quantify the abundance of differentially expressed genes from **NOT normalized** reads and aligned with **Bowtie 1**.

```
read.table("./data/summary.br.nonorm.44062.txt") %>%
```

```

ggplot(aes(
  x = V1,
  y = V10,
  fill = V7)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V1,
    y = V10,
    ymax = V10,
    label = V10,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V4 + V5 + V6 + V8,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



53
54 The R package *eXpress* is solely used to quantify the abundance of differentially expressed genes from
55 **normalized** reads and aligned with **Bowtie 1**.

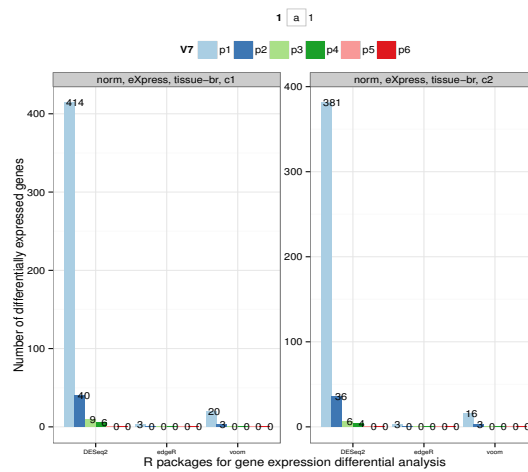
```
read.table("./data/summary.br.norm.44060.txt") %>%
```



```

ggplot(aes(
  x = V1,
  y = V10,
  fill = V7)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V1,
    y = V10,
    ymax = V10,
    label = V10,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V4 + V5 + V6 + V8,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



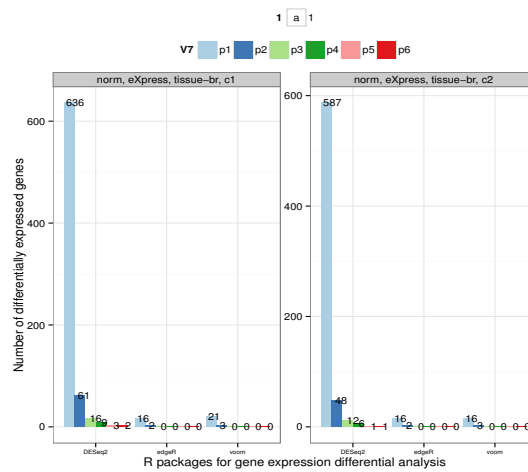
56 The R package *eXpress* is solely used to quantify the abundance of differentially expressed genes from
 57 **normalized** reads and aligned with **Bowtie 2**.
 58

```
read.table("./data/summary.br.norm.44061.txt") %>%
```

```

ggplot(aes(
  x = V1,
  y = V10,
  fill = V7)) +
  geom_bar(stat = "identity",
    position = "dodge") +
  geom_text(aes(x = V1,
    y = V10,
    ymax = V10,
    label = V10,
    size = 1,
    hjust = 0),
    position = position_dodge(width = 1)) +
  facet_wrap(~ V4 + V5 + V6 + V8,
    ncol = 2,
    scales = "free") +
  scale_fill_brewer(type = "qual", palette = "Paired") +
  theme_bw() +
  theme(legend.position = "top",
    axis.text.x = element_text(vjust = .5,
      size = 6)) +
  labs(x = "R packages for gene expression differential analysis",
    y = "Number of differentially expressed genes")

```



3 Identifying foreign transcripts in oyster cells

Microscopy protocols identified the presence of a form of parasitic worm in healthy oyster tissues. As a result we might have a parasitic contamination of our RNA-seq reads. Fortunately, this can give a chance to extract the foreign RNA reads, which are already sequenced, assemble them so we identify later on the parasitic species.

We gathered 5 genomes of parasitic worms [first link](#), [second link](#). We used 2 separate strategies to assemble contigs specific for each genome, i- using genome-guided assembly by mapping reads to genome or ii- mapping contigs directly to genome.

3.1 Genome-guided assembly

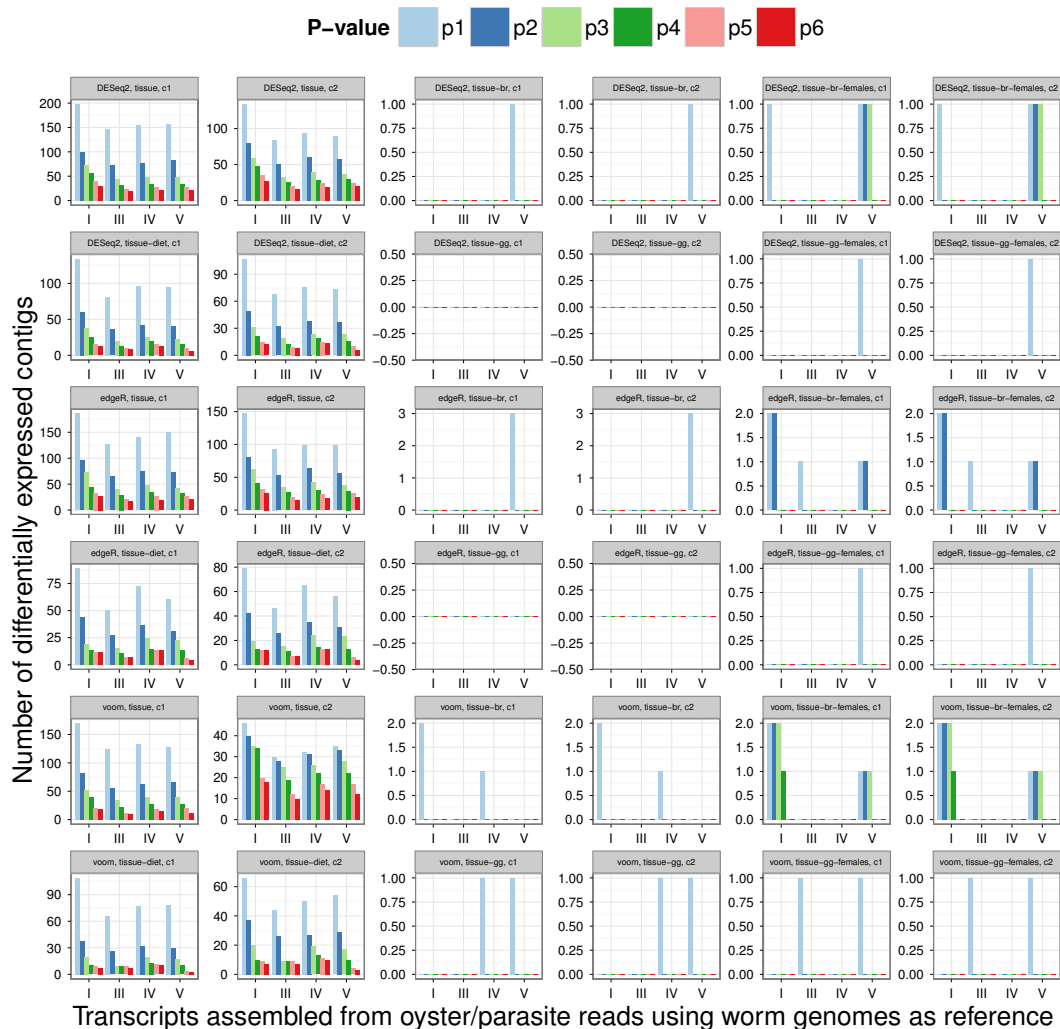
We used bwa to map reads to each genome separately, then use only the proper mapped mates to assemble a special transcriptome using trinity. Those selected reads were then aligned to each of the five new transcriptomes separately. Abundance of reads and differentially expressed transcripts were finally identified.

1. *Clonorchis sinensis*
2. *Opisthorchis viverrini*
3. *Schistosoma haematobium*
4. *Schistosoma japonicum*
5. *Schistosoma mansoni*

```

p1 <- read.table("./data/summary.eXpress.134221.txt")
p3 <- read.table("./data/summary.eXpress.134239.txt")
p4 <- read.table("./data/summary.eXpress.134248.txt")
p5 <- read.table("./data/summary.eXpress.135450.txt")
p1$V9 <- c("I")
p3$V9 <- c("III")
p4$V9 <- c("IV")
p5$V9 <- c("V")
rbind(p1,p3,p4,p5) %>%
  ggplot(aes(x = V9,
             y = V7,
             fill = V4)) +
  theme_bw() +
  geom_bar(stat = "identity",
           position = "dodge") +
  facet_wrap(~ V1 + V3 + V5,
             ncol = 6,
             scales = "free") +
  labs(x = "Transcripts assembled from oyster/parasite reads using worm genomes as reference",
       y = "Number of differentially expressed contigs") +
  theme(legend.position = "top",
        axis.text.x = element_text(vjust = .5, size = 6),
        axis.text.y = element_text(vjust = .5, size = 6),
        strip.text = element_text(size = 4),
        axis.ticks = element_line(size = .2),
        axis.ticks.length = unit(.05, "cm")) +
  scale_fill_brewer(type = "qual", palette = "Paired",
                    name = "P-value")

```



3.2 Aligning contigs to worm genomes

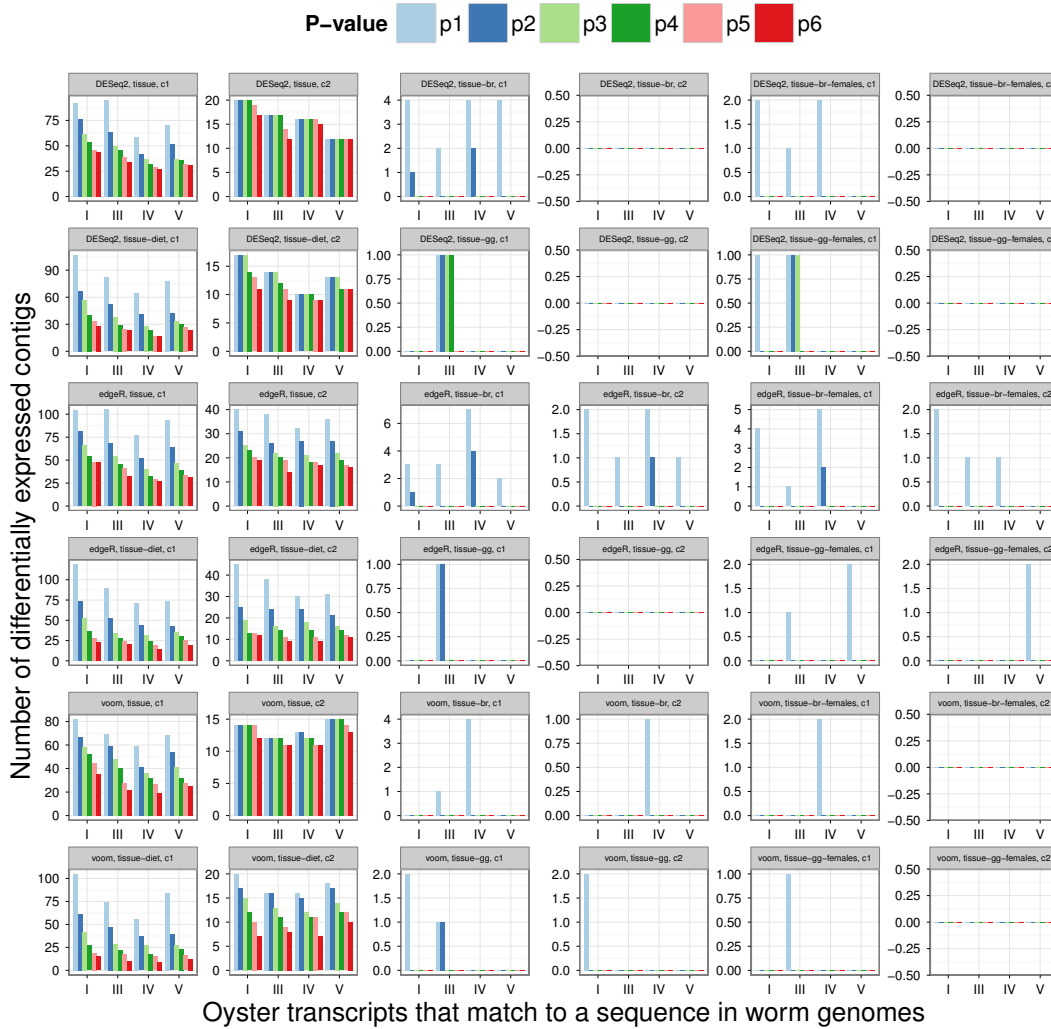
Contigs were generated by assembling raw sequencing reads from all tissues and dietary conditions. These transcripts were then aligned to each worm genome separately. We chose the alignments that match at least 500 sequence length.

```
p1 <- read.table("../data/summary.eXpress.134745.txt")
```

```

p3 <- read.table("./data/summary.eXpress.134744.txt")
p4 <- read.table("./data/summary.eXpress.135311.txt")
p5 <- read.table("./data/summary.eXpress.135424.txt")
p1$V9 <- c("I")
p3$V9 <- c("III")
p4$V9 <- c("IV")
p5$V9 <- c("V")
rbind(p1,p3,p4,p5) %>%
  ggplot(aes(x = V9,
             y = V7,
             fill = V4)) +
  theme_bw() +
  geom_bar(stat = "identity",
           position = "dodge") +
  facet_wrap(~ V1 + V3 + V5,
            ncol = 6,
            scales = "free") +
  labs(x = "Oyster transcripts that match to a sequence in worm genomes",
       y = "Number of differentially expressed contigs") +
  theme(legend.position = "top",
        axis.text.x = element_text(vjust = .5, size = 6),
        axis.text.y = element_text(vjust = .5, size = 6),
        strip.text = element_text(size = 4),
        axis.ticks = element_line(size = .2),
        axis.ticks.length = unit(.05, "cm")) +
  scale_fill_brewer(type = "qual", palette = "Paired",
                   name = "P-value")

```

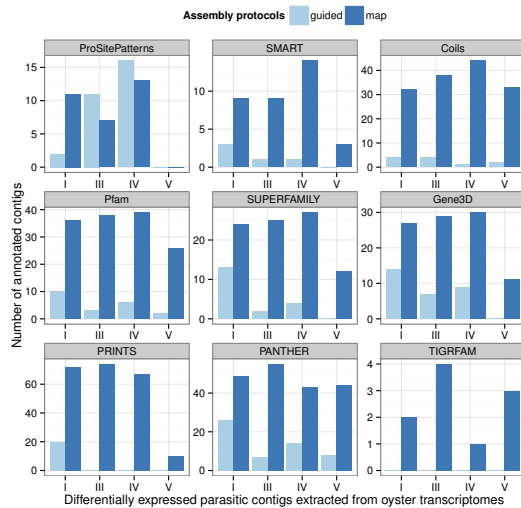


84

85 3.3 Annotating genes found from mapping parasitic genomes to oyster reads

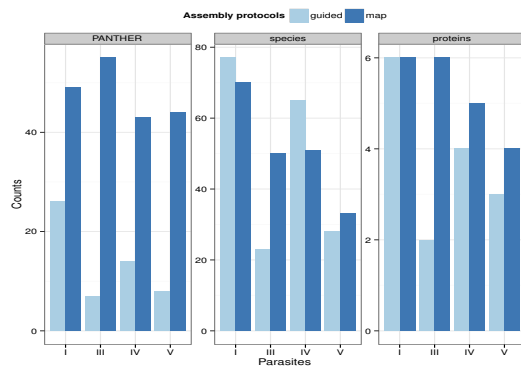
86 We used protein databases from ten different sources: **ProDom**, **PANTHER**, **TIGRFAM**, **SUPERFAMILY**,
 87 **PIRSF**, **Gene3D**, **Pfam**, **SMART**, **PROSITEPROFILES**, **PROSITEPATTERNS**, **COILS**. Hidden Markov
 88 models (HMM) were used to find annotations belonging to differentially expressed genes in section 3.1
 89 and 3.2. We chose tissue specific differentially expressed genes at a p-value of 10e-3 and a c-fold of 2.
 90 DESeq2 R packages generated these gene selections between ganglia and gill tissues.

```
read.table("./data/parasite-ips.txt", header = T) %>%
  gather("database", "counts", 2:10) %>%
  ggplot(aes(x = parasite,
             y = counts,
             fill = assembly)) +
  theme_bw() +
  geom_bar(stat = "identity",
           position = "dodge") +
  facet_wrap(~ database,
             ncol = 3,
             scales = "free") +
  labs(x = "Differentially expressed parasitic contigs extracted from oyster transcriptomes",
       y = "Number of annotated contigs") +
  theme(legend.position = "top") +
  scale_fill_brewer(type = "qual", palette = "Paired",
                   name = "Assembly protocols")
```



The Panther database has an extended collection of species, functional domains, and GO-terms compared to other databases. Here is a short summary of what was found in Panther at e-value of 10e-10 and minimum amino acid alignment length of 20.

```
read.table("./data/parasite-panther.txt", header = T) %>%
  gather("category", "count", 2:4) %>%
  ggplot(aes(x = parasite,
             y = count,
             fill = assembly)) +
  theme_bw() +
  geom_bar(stat = "identity",
           position = "dodge") +
  facet_wrap(~ category,
             ncol = 3,
             scales = "free") +
  labs(x = "Parasites",
       y = "Counts") +
  theme(legend.position = "top") +
  scale_fill_brewer(type = "qual", palette = "Paired",
                   name = "Assembly protocols")
```



3.4 How many reads map by pair mates to each parasite genome?

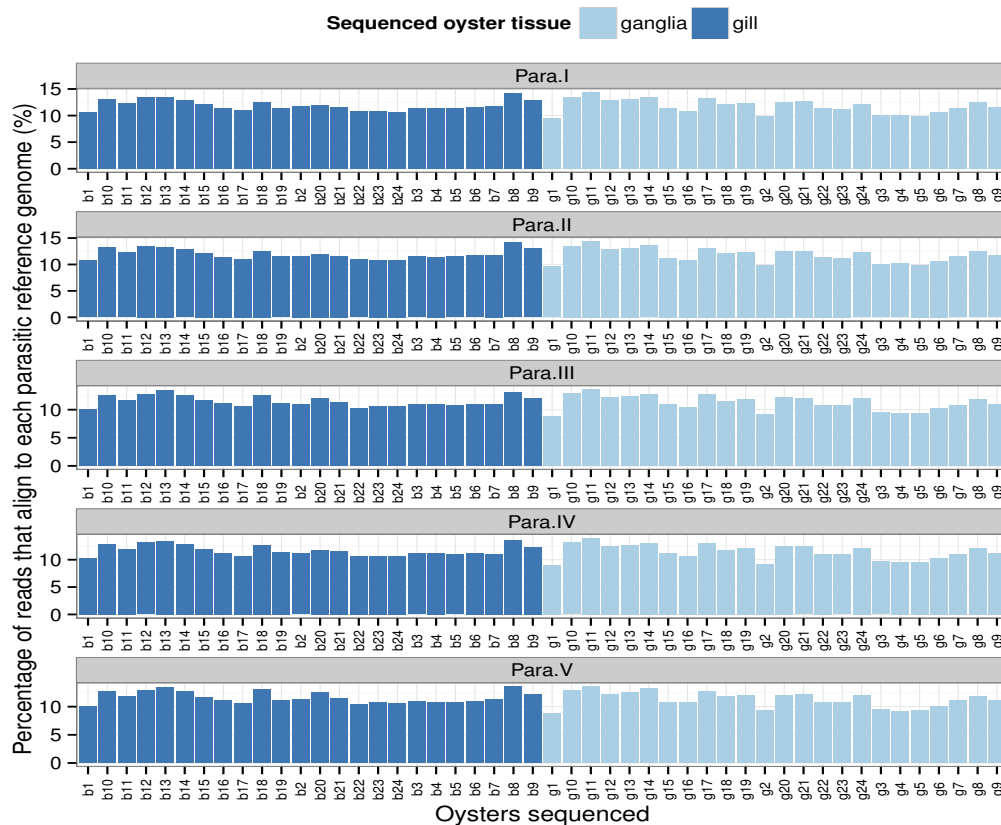
Only pair reads where both mates map in a forward and reverse fashion are displayed.

```
df <- read.table("./data/parasite.mates.mapping.txt", header = T)
```

```

dfx <- sapply(df[,c(2:6)], function(x) {(x/df$Total)*100})
dfx <- data.frame(dfx, df$Tissue, df$Sample)
#require(scales)
#dfy[order(dfy$counts, decreasing=FALSE),] %>% head
gather(dfx, "parasite", "counts", 1:5) %>%
#arrange(desc(counts)) %>%
  ggplot(aes(x = df.Sample,
             y = counts,
             fill = df.Tissue)) +
  theme_bw() +
  geom_bar(stat = "identity",
           position = "dodge") +
  facet_wrap(~ parasite,
             ncol = 1,
             scales = "free") +
  labs(x = "Oysters sequenced",
       y = "Percentage of reads that align to each parasitic reference genome (%)") +
  theme(legend.position = "top",
        axis.text.x = element_text(angle = 90,
                                     vjust = .5, size = 7)) +
  scale_fill_brewer(type = "qual", palette = "Paired",
                    name = "Sequenced oyster tissue")
# scale_y_continuous(labels = scales::percent)

```



98

99 4 Network inference

100 We will be using weighted genes co-expression analyses from Zhang and Horraht 2005. Machine learning
 101 regularization techniques will be applied after to emphasize selection of significant gene modules.

102 4.1 Correlation analysis

103 The analysis sill include: Sample verification and clustering. Similarity matrix between samples. Conver-
 104 tion of similarity matrix to adjacency matrix. Detection of co-expression modules.

```

counts <- read.table("./data/diffExpr.P1e-4_C2.matrix.log2.dat", header = T)
heatmap.2(cor(counts), trace = 'none')

```


sessionInfo()

R version 3.2.1 (2015-06-18)

Platform: x86_64-unknown-linux-gnu (64-bit)

Running under: elementary OS Luna

locale:

[1] LC_CTYPE=en_US.UTF-8	LC_NUMERIC=C
[3] LC_TIME=en_US.UTF-8	LC_COLLATE=en_US.UTF-8
[5] LC_MONETARY=en_US.UTF-8	LC_MESSAGES=en_US.UTF-8
[7] LC_PAPER=en_US.UTF-8	LC_NAME=en_US.UTF-8
[9] LC_ADDRESS=en_US.UTF-8	LC_TELEPHONE=en_US.UTF-8
[11] LC_MEASUREMENT=en_US.UTF-8	LC_IDENTIFICATION=en_US.UTF-8

attached base packages:

[1] grid	stats	graphics	grDevices	utils	datasets
[7] methods	base				

other attached packages:

[1] stringr_1.0.0	reshape2_1.4.1	WGCNA_1.47
[4] RSQLite_1.0.0	DBI_0.3.1	fastcluster_1.1.16
[7] dynamicTreeCut_1.62	gplots_2.17.0	tidyr_0.2.0
[10] dplyr_0.4.2	latticeExtra_0.6-26	RColorBrewer_1.1-2
[13] glmnet_2.0-2	foreach_1.4.2	Matrix_1.2-1
[16] leaps_2.9	caret_6.0-47	ggplot2_1.0.1
[19] lattice_0.20-31	xlsx_0.5.7	xlsxjars_0.6.1
[22] rJava_0.9-6	knitr_1.10.5	RevoUtilsMath_3.2.1

loaded via a namespace (and not attached):

[1] Biobase_2.30.0	splines_3.2.1	gtools_3.5.0
[4] Formula_1.2-1	assertthat_0.1	highr_0.5
[7] stats4_3.2.1	impute_1.44.0	quantreg_5.11
[10] digest_0.6.8	minqa_1.2.4	colorspace_1.2-6
[13] preprocessCore_1.32.0	plyr_1.8.3	BradleyTerry2_1.0-6
[16] SparseM_1.6	GO.db_3.2.2	scales_0.2.5
[19] gdata_2.16.1	brglm_0.5-9	lme4_1.1-8
[22] mgcv_1.8-6	IRanges_2.4.8	car_2.0-25
[25] lazyeval_0.1.10	nnet_7.3-10	BiocGenerics_0.16.1
[28] pbkrtest_0.4-2	proto_0.3-10	survival_2.38-2
[31] magrittr_1.5	evaluate_0.7	doParallel_1.0.8
[34] nlme_3.1-121	MASS_7.3-41	foreign_0.8-64
[37] tools_3.2.1	formatR_1.2	matrixStats_0.14.2
[40] S4Vectors_0.8.11	munsell_0.4.2	cluster_2.0.2
[43] AnnotationDbi_1.32.3	caTools_1.17.1	nloptr_1.0.4
[46] iterators_1.0.7	labeling_0.3	bitops_1.0-6
[49] gtable_0.1.2	codetools_0.2-11	reshape_0.8.5
[52] R6_2.0.1	gridExtra_0.9.1	Hmisc_3.16-0
[55] KernSmooth_2.23-15	stringi_0.5-5	parallel_3.2.1
[58] Rcpp_0.11.6	rpart_4.1-10	acepack_1.3-3.3