

Class 14 RNA Mini-project

Zainab Fatima (PID: A16880407)

2025-02-20

Table of contents

Background	1
Data Import	2
Inspect and tidy	2
Setup for DESeq	3
Run DESeq	3
Volcano plot of results	5
Gene annotations	7
Pathway Analysis	8
Gene Ontology (GO)	20
Reactome Analysis	21
GO online (optional)	22

Background

The data for for hands-on session comes from GEO entry: GSE37704, which is associated with the following publication:

Trapnell C, Hendrickson DG, Sauvageau M, Goff L et al. “Differential analysis of gene regulation at transcript resolution with RNA-seq”. Nat Biotechnol 2013 Jan;31(1):46-53. PMID: 23222703

The authors report on differential analysis of lung fibroblasts in response to loss of the developmental transcription factor HOXA1. Their results and others indicate that HOXA1 is required for lung fibroblast and HeLa cell cycle progression. In particular their analysis show that “loss of HOXA1 results in significant expression level changes in thousands of individual transcripts, along with isoform switching events in key regulators of the cell cycle”. For our session we have used their Sailfish gene-level estimated counts and hence are restricted to protein-coding genes only.

Data Import

```
counts <- read.csv("GSE37704_featurecounts.csv", row.names = 1, stringsAsFactors = F, header = 1)
colData <- read.csv("GSE37704_metadata.csv", stringsAsFactors = F, header = T)
```

Inspect and tidy

Q. Complete the code below to remove the troublesome first column from countData

```
colnames(counts)
```

```
[1] "length"      "SRR493366" "SRR493367" "SRR493368" "SRR493369" "SRR493370"
[7] "SRR493371"
```

```
#need to remove length column
```

```
countData <- counts[,-1]
```

```
head(countData)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000186092	0	0	0	0	0	0
ENSG00000279928	0	0	0	0	0	0
ENSG00000279457	23	28	29	29	28	46
ENSG00000278566	0	0	0	0	0	0
ENSG00000273547	0	0	0	0	0	0
ENSG00000187634	124	123	205	207	212	258

Check for matching countData and coldata

```
colnames(countData) %in% colData$id
```

```
[1] TRUE TRUE TRUE TRUE TRUE TRUE
```

```
colData[match(colnames(countData),colData$id),]
```

	id	condition
1	SRR493366	control_sirna
2	SRR493367	control_sirna
3	SRR493368	control_sirna
4	SRR493369	hoxa1_kd
5	SRR493370	hoxa1_kd
6	SRR493371	hoxa1_kd

Q. Complete the code below to filter countData to exclude genes (i.e. rows) where we have 0 read count across all samples (i.e. columns).

```
to.keep.inds <- rowSums(countData) > 0
```

Q. Complete the code below to filter countData to exclude genes (i.e. rows) where we have 0 read count across all samples (i.e. columns).

```
new.counts <- countData[to.keep.inds, ]
head(new.counts)
```

	SRR493366	SRR493367	SRR493368	SRR493369	SRR493370	SRR493371
ENSG00000279457	23	28	29	29	28	46
ENSG00000187634	124	123	205	207	212	258
ENSG00000188976	1637	1831	2383	1226	1326	1504
ENSG00000187961	120	153	180	236	255	357
ENSG00000187583	24	48	65	44	48	64
ENSG00000187642	4	9	16	14	16	16

```
nrow(new.counts)
```

```
[1] 15975
```

Setup for DESeq

```
library(DESeq2)
```

Run DESeq

```
dds <- DESeqDataSetFromMatrix(countData=countData,  
                              colData=colData,  
                              design=~condition)
```

Warning in DESeqDataSet(se, design = design, ignoreRank): some variables in design formula are characters, converting to factors

```
dds <- DESeq(dds)
```

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

```
dds
```

```
class: DESeqDataSet  
dim: 19808 6  
metadata(1): version  
assays(4): counts mu H cooks  
rownames(19808): ENSG00000186092 ENSG00000279928 ... ENSG00000277475  
               ENSG00000268674  
rowData names(22): baseMean baseVar ... deviance maxCooks  
colnames(6): SRR493366 SRR493367 ... SRR493370 SRR493371  
colData names(3): id condition sizeFactor
```

```
res <- results(dds, contrast=c("condition", "hoxa1_kd", "control_sirna"))
```

```
head(res)
```

```
log2 fold change (MLE): condition hoxa1_kd vs control_sirna
```

```
Wald test p-value: condition hoxa1 kd vs control sirna
```

```
DataFrame with 6 rows and 6 columns
```

	baseMean	log2FoldChange	lfcSE	stat	pvalue
	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>
ENSG00000186092	0.0000	NA	NA	NA	NA
ENSG00000279928	0.0000	NA	NA	NA	NA
ENSG00000279457	29.9136	0.179257	0.324822	0.551863	0.58104205
ENSG00000278566	0.0000	NA	NA	NA	NA
ENSG00000273547	0.0000	NA	NA	NA	NA
ENSG00000187634	183.2296	0.426457	0.140266	3.040350	0.00236304
	padj				
	<numeric>				
ENSG00000186092	NA				
ENSG00000279928	NA				
ENSG00000279457	0.68707978				
ENSG00000278566	NA				
ENSG00000273547	NA				
ENSG00000187634	0.00516278				

Q. Call the `summary()` function on your results to get a sense of how many genes are up or down-regulated at the default 0.1 p-value cutoff.

```
summary(res)
```

```
out of 15975 with nonzero total read count
```

```
adjusted p-value < 0.1
```

```
LFC > 0 (up)      : 4349, 27%
```

```
LFC < 0 (down)    : 4393, 27%
```

```
outliers [1]      : 0, 0%
```

```
low counts [2]    : 1221, 7.6%
```

```
(mean count < 0)
```

```
[1] see 'cooksCutoff' argument of ?results
```

```
[2] see 'independentFiltering' argument of ?results
```

Volcano plot of results

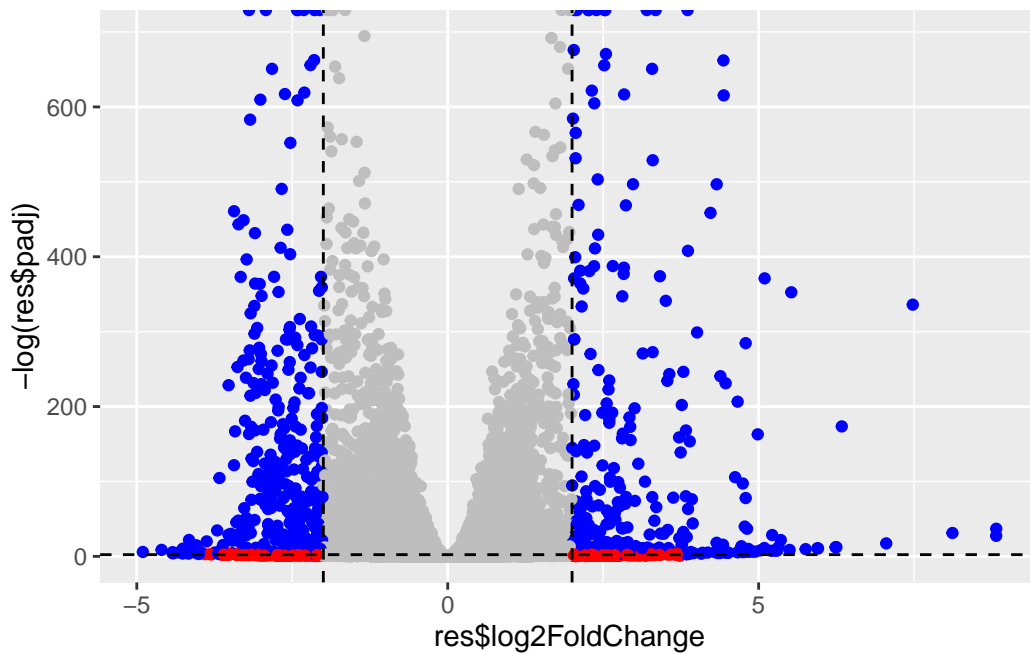
Q. Improve this plot by completing the below code, which adds color and axis labels

```
library(ggplot2)
```

```
mycols <- rep("gray", nrow(res))  
#if my log2fold change >2 --> change it to red  
mycols[ abs(res$log2FoldChange) > 2 ] <- "red"  
inds <- (res$padj < 0.05) & (abs(res$log2FoldChange) > 2 )  
mycols[ inds ] <- "blue"
```

```
ggplot(res) +  
  aes(x = res$log2FoldChange, y = -log(res$padj)) +  
  geom_point(col = mycols) +  
  geom_vline(xintercept = 2, linetype = "dashed") +  
  geom_vline(xintercept = -2, linetype = "dashed") +  
  geom_hline(yintercept = -log(0.1), linetype = "dashed")
```

Warning: Removed 5054 rows containing missing values or values outside the scale range (`geom_point()`).



Gene annotations

Q. Use the `mapIds()` function multiple times to add SYMBOL, ENTREZID and GENENAME annotation to our results by completing the code below.

```
library(AnnotationDbi)
library(org.Hs.eg.db)
```

```
res$symbol <- mapIds(org.Hs.eg.db,
  keys = rownames(res),
  keytype = "ENSEMBL",
  column = "SYMBOL")
```

'select()' returned 1:many mapping between keys and columns

```
res$genename <- mapIds(org.Hs.eg.db,
  keys = rownames(res),
  keytype = "ENSEMBL",
  column = "GENENAME")
```

'select()' returned 1:many mapping between keys and columns

```
res$entrez <- mapIds(org.Hs.eg.db,
  keys = rownames(res),
  keytype = "ENSEMBL",
  column = "ENTREZID")
```

'select()' returned 1:many mapping between keys and columns

```
head(res)
```

log2 fold change (MLE): condition hoxa1_kd vs control_sirna

Wald test p-value: condition hoxa1 kd vs control sirna

DataFrame with 6 rows and 9 columns

	baseMean	log2FoldChange	lfcSE	stat	pvalue
	<numeric>	<numeric>	<numeric>	<numeric>	<numeric>
ENSG00000186092	0.0000	NA	NA	NA	NA

ENSG00000279928	0.0000	NA	NA	NA	NA
ENSG00000279457	29.9136	0.179257	0.324822	0.551863	0.58104205
ENSG00000278566	0.0000	NA	NA	NA	NA
ENSG00000273547	0.0000	NA	NA	NA	NA
ENSG00000187634	183.2296	0.426457	0.140266	3.040350	0.00236304
	padj	symbol		genename	entrez
	<numeric>	<character>		<character>	<character>
ENSG00000186092	NA	OR4F5	olfactory receptor f..		79501
ENSG00000279928	NA	NA		NA	NA
ENSG00000279457	0.68707978	NA		NA	NA
ENSG00000278566	NA	NA		NA	NA
ENSG00000273547	NA	NA		NA	NA
ENSG00000187634	0.00516278	SAMD11	sterile alpha motif ..		148398

Q. Finally for this section let's reorder these results by adjusted p-value and save them to a CSV file in your current project directory.

```
res = res[order(res$pvalue),]
write.csv(res, file="deseq_results.csv")
```

Pathway Analysis

```
library(pathview)
```

```
#####
Pathview is an open source software package distributed under GNU General
Public License version 3 (GPLv3). Details of GPLv3 is available at
http://www.gnu.org/licenses/gpl-3.0.html. Particullary, users are required to
formally cite the original Pathview paper (not just mention it) in publications
or products. For details, do citation("pathview") within R.
```

The pathview downloads and uses KEGG data. Non-academic uses may require a KEGG license agreement (details at <http://www.kegg.jp/kegg/legal.html>).

```
#####
```

```
library(gage)
```



```
library(gageData)
```

```
data(kegg.sets.hs)  
data(sigmet.idx.hs)
```

```
# Focus on signaling and metabolic pathways only  
kegg.sets.hs = kegg.sets.hs[sigmet.idx.hs]
```

```
# Examine the first 3 pathways  
head(kegg.sets.hs, 3)
```

```
$`hsa00232 Caffeine metabolism`
```

```
[1] "10" "1544" "1548" "1549" "1553" "7498" "9"
```

```
$`hsa00983 Drug metabolism - other enzymes`
```

```
[1] "10" "1066" "10720" "10941" "151531" "1548" "1549" "1551"  
[9] "1553" "1576" "1577" "1806" "1807" "1890" "221223" "2990"  
[17] "3251" "3614" "3615" "3704" "51733" "54490" "54575" "54576"  
[25] "54577" "54578" "54579" "54600" "54657" "54658" "54659" "54963"  
[33] "574537" "64816" "7083" "7084" "7172" "7363" "7364" "7365"  
[41] "7366" "7367" "7371" "7372" "7378" "7498" "79799" "83549"  
[49] "8824" "8833" "9" "978"
```

```
$`hsa00230 Purine metabolism`
```

```
[1] "100" "10201" "10606" "10621" "10622" "10623" "107" "10714"  
[9] "108" "10846" "109" "111" "11128" "11164" "112" "113"  
[17] "114" "115" "122481" "122622" "124583" "132" "158" "159"  
[25] "1633" "171568" "1716" "196883" "203" "204" "205" "221823"  
[33] "2272" "22978" "23649" "246721" "25885" "2618" "26289" "270"  
[41] "271" "27115" "272" "2766" "2977" "2982" "2983" "2984"  
[49] "2986" "2987" "29922" "3000" "30833" "30834" "318" "3251"  
[57] "353" "3614" "3615" "3704" "377841" "471" "4830" "4831"  
[65] "4832" "4833" "4860" "4881" "4882" "4907" "50484" "50940"  
[73] "51082" "51251" "51292" "5136" "5137" "5138" "5139" "5140"  
[81] "5141" "5142" "5143" "5144" "5145" "5146" "5147" "5148"  
[89] "5149" "5150" "5151" "5152" "5153" "5158" "5167" "5169"  
[97] "51728" "5198" "5236" "5313" "5315" "53343" "54107" "5422"  
[105] "5424" "5425" "5426" "5427" "5430" "5431" "5432" "5433"  
[113] "5434" "5435" "5436" "5437" "5438" "5439" "5440" "5441"  
[121] "5471" "548644" "55276" "5557" "5558" "55703" "55811" "55821"  
[129] "5631" "5634" "56655" "56953" "56985" "57804" "58497" "6240"
```

```
[137] "6241"    "64425"   "646625"  "654364"  "661"     "7498"    "8382"    "84172"
[145] "84265"   "84284"   "84618"   "8622"     "8654"    "87178"   "8833"    "9060"
[153] "9061"    "93034"   "953"      "9533"     "954"     "955"     "956"     "957"
[161] "9583"    "9615"
```

The main `gage()` function requires a named vector of fold changes, where the names of the values are the Entrez gene IDs.

```
foldchanges = res$log2FoldChange
names(foldchanges) = res$entrez
head(foldchanges)
```

```
      1266      54855      1465      51232      2034      2317
-2.422719  3.201955 -2.313738 -2.059631 -1.888019 -1.649792
```

Now, let's run the `gage` pathway analysis.

```
#Get the results

keggres = gage(foldchanges, gsets=kegg.sets.hs)
```

```
attributes(keggres)
```

```
$names
[1] "greater" "less"    "stats"
```

```
# Look at the first few down (less) pathways
head(keggres$less)
```

	p.geomean	stat.mean	p.val
hsa04110 Cell cycle	7.077982e-06	-4.432593	7.077982e-06
hsa03030 DNA replication	9.424076e-05	-3.951803	9.424076e-05
hsa03013 RNA transport	1.160132e-03	-3.080629	1.160132e-03
hsa04114 Oocyte meiosis	2.563806e-03	-2.827297	2.563806e-03
hsa03440 Homologous recombination	3.066756e-03	-2.852899	3.066756e-03
hsa00010 Glycolysis / Gluconeogenesis	4.360092e-03	-2.663825	4.360092e-03
	q.val	set.size	exp1
hsa04110 Cell cycle	0.001160789	124	7.077982e-06
hsa03030 DNA replication	0.007727742	36	9.424076e-05

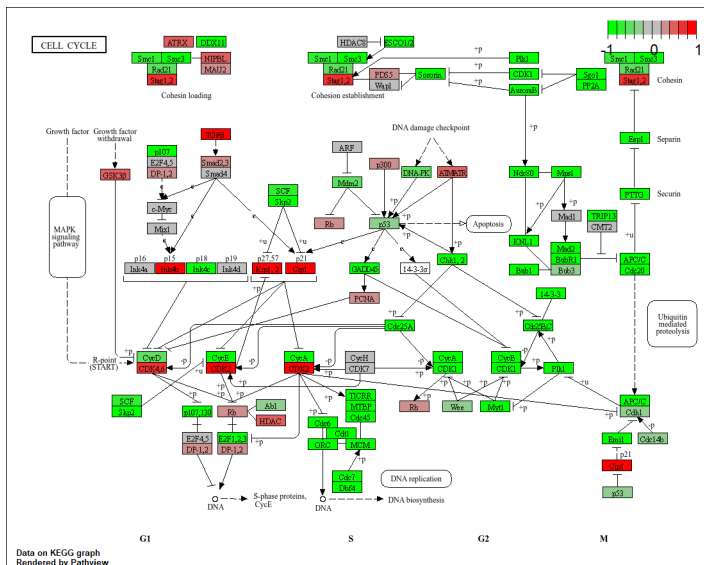
hsa03013	RNA transport	0.063420543	149	1.160132e-03
hsa04114	Oocyte meiosis	0.100589607	112	2.563806e-03
hsa03440	Homologous recombination	0.100589607	28	3.066756e-03
hsa00010	Glycolysis / Gluconeogenesis	0.119175854	65	4.360092e-03

```
pathview(gene.data=foldchanges, pathway.id="hsa04110")
```

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04110.pathview.png



We can play with the other input arguments to `pathview()` to change the display in various ways including generating a PDF graph. For example:

```
pathview(gene.data=foldchanges, pathway.id="hsa04110", kegg.native=FALSE)
```

'select()' returned 1:1 mapping between keys and columns

Warning: reconcile groups sharing member nodes!

```
      [,1] [,2]
[1,] "9"  "300"
[2,] "9"  "306"
```

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04110.pathview.pdf

Now, let's pull up the top 5 upregulated pathways.

```
keggrespathways <- rownames(keggres$greater)[1:5]
```

```
# Extract the 8 character long IDs part of each string
keggresids = substr(keggrespathways, start=1, stop=8)
keggresids
```

```
[1] "hsa04740" "hsa04640" "hsa00140" "hsa04630" "hsa04976"
```

```
pathview(gene.data=foldchanges, pathway.id=keggresids, species="hsa")
```

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04740.pathview.png

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04640.pathview.png

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa00140.pathview.png

'select()' returned 1:1 mapping between keys and columns

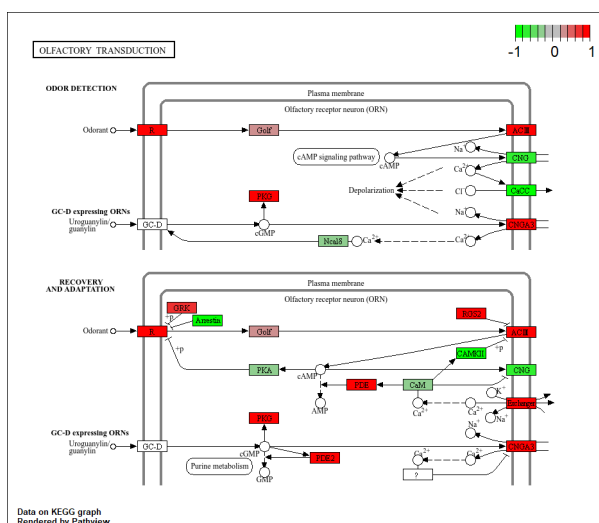
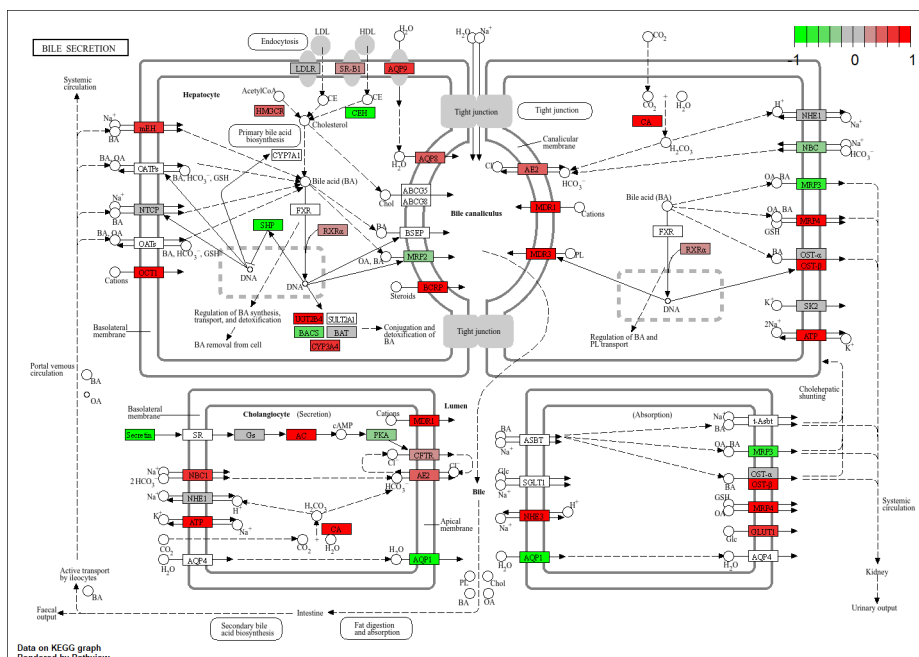
Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04630.pathview.png

'select()' returned 1:1 mapping between keys and columns

Info: Working in directory C:/BIMM 143/Class 14

Info: Writing image file hsa04976.pathview.png



Q. Can you do the same procedure as above to plot the pathview figures for the top 5 down-regulated pathways?

```
keggrespathways2 <- rownames(keggres$less)[1:5]
```

```
keggresids2 = substr(keggrespathways2, start=1, stop=8)
keggresids2
```



```
[1] "hsa04110" "hsa03030" "hsa03013" "hsa04114" "hsa03440"
```

```
pathview(gene.data=foldchanges, pathway.id=keggresids2, species="hsa")
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory C:/BIMM 143/Class 14
```

```
Info: Writing image file hsa04110.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory C:/BIMM 143/Class 14
```

```
Info: Writing image file hsa03030.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory C:/BIMM 143/Class 14
```

```
Info: Writing image file hsa03013.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

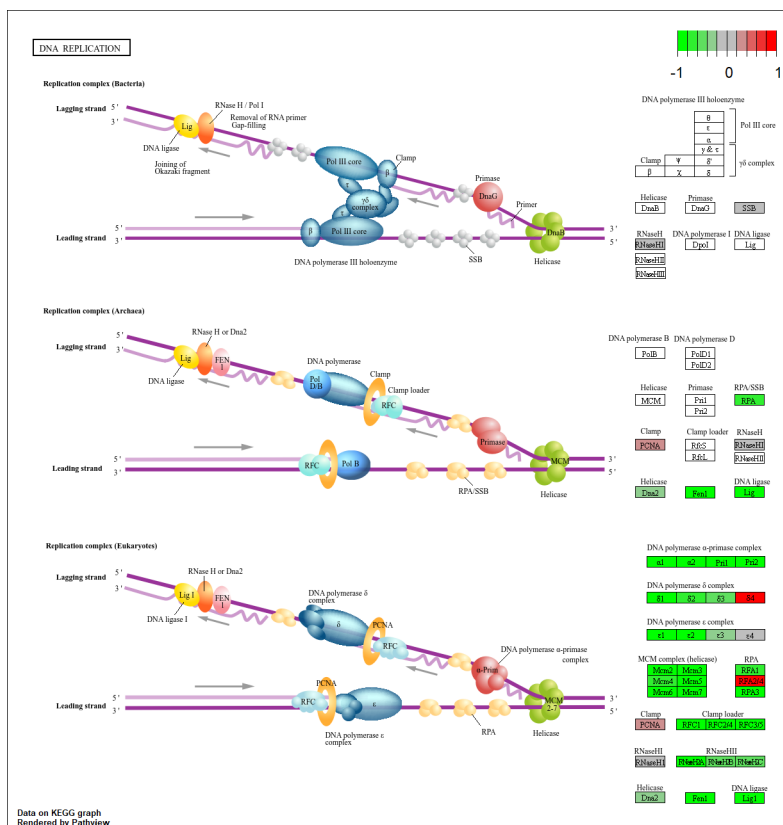
```
Info: Working in directory C:/BIMM 143/Class 14
```

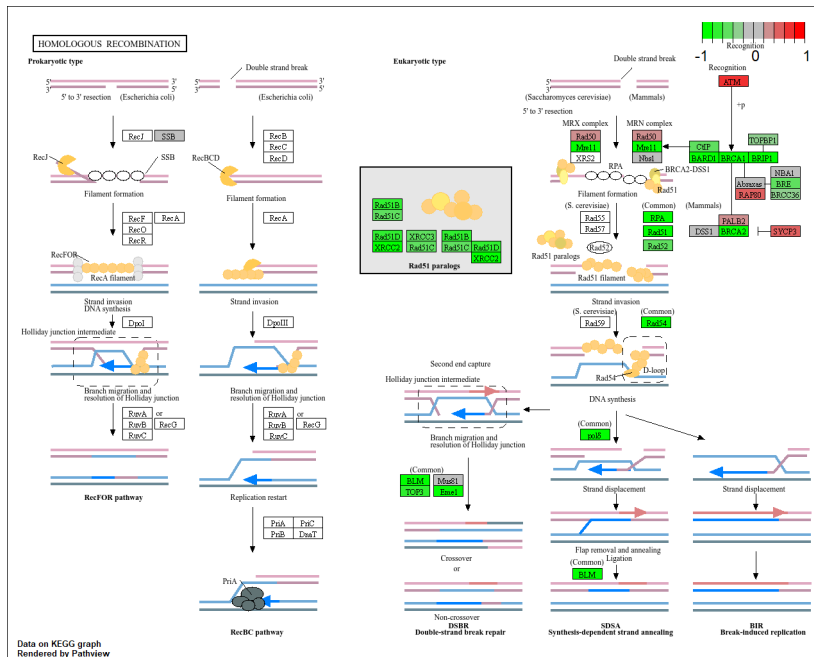
```
Info: Writing image file hsa04114.pathview.png
```

```
'select()' returned 1:1 mapping between keys and columns
```

```
Info: Working in directory C:/BIMM 143/Class 14
```

```
Info: Writing image file hsa03440.pathview.png
```





Gene Ontology (GO)

```
data(go.sets.hs)
data(go.subs.hs)

# Focus on Biological Process subset of GO
gobpsets = go.sets.hs[go.subs.hs$BP]

gobpres = gage(foldchanges, gsets=gobpsets, same.dir=TRUE)

lapply(gobpres, head)
```

\$greater

		p.geomean	stat.mean	p.val
G0:0007156	homophilic cell adhesion	1.734864e-05	4.210777	1.734864e-05
G0:0048729	tissue morphogenesis	5.407952e-05	3.888470	5.407952e-05
G0:0002009	morphogenesis of an epithelium	5.727599e-05	3.878706	5.727599e-05
G0:0030855	epithelial cell differentiation	2.053700e-04	3.554776	2.053700e-04
G0:0060562	epithelial tube morphogenesis	2.927804e-04	3.458463	2.927804e-04
G0:0048598	embryonic morphogenesis	2.959270e-04	3.446527	2.959270e-04
		q.val	set.size	exp1

G0:0007156	homophilic cell adhesion	0.07584825	137	1.734864e-05
G0:0048729	tissue morphogenesis	0.08347021	483	5.407952e-05
G0:0002009	morphogenesis of an epithelium	0.08347021	382	5.727599e-05
G0:0030855	epithelial cell differentiation	0.16449701	299	2.053700e-04
G0:0060562	epithelial tube morphogenesis	0.16449701	289	2.927804e-04
G0:0048598	embryonic morphogenesis	0.16449701	498	2.959270e-04

\$less

		p.geomean	stat.mean	p.val
G0:0048285	organelle fission	6.626774e-16	-8.170439	6.626774e-16
G0:0000280	nuclear division	1.797050e-15	-8.051200	1.797050e-15
G0:0007067	mitosis	1.797050e-15	-8.051200	1.797050e-15
G0:0000087	M phase of mitotic cell cycle	4.757263e-15	-7.915080	4.757263e-15
G0:0007059	chromosome segregation	1.081862e-11	-6.974546	1.081862e-11
G0:0051301	cell division	8.718528e-11	-6.455491	8.718528e-11

		q.val	set.size	exp1
G0:0048285	organelle fission	2.618901e-12	386	6.626774e-16
G0:0000280	nuclear division	2.618901e-12	362	1.797050e-15
G0:0007067	mitosis	2.618901e-12	362	1.797050e-15
G0:0000087	M phase of mitotic cell cycle	5.199689e-12	373	4.757263e-15
G0:0007059	chromosome segregation	9.459800e-09	146	1.081862e-11
G0:0051301	cell division	6.352901e-08	479	8.718528e-11

\$stats

		stat.mean	exp1
G0:0007156	homophilic cell adhesion	4.210777	4.210777
G0:0048729	tissue morphogenesis	3.888470	3.888470
G0:0002009	morphogenesis of an epithelium	3.878706	3.878706
G0:0030855	epithelial cell differentiation	3.554776	3.554776
G0:0060562	epithelial tube morphogenesis	3.458463	3.458463
G0:0048598	embryonic morphogenesis	3.446527	3.446527

Reactome Analysis

```
sig_genes <- res[res$padj <= 0.05 & !is.na(res$padj), "symbol"]
print(paste("Total number of significant genes:", length(sig_genes)))
```

```
[1] "Total number of significant genes: 8146"
```

```
write.table(sig_genes, file="significant_genes.txt", row.names=FALSE, col.names=FALSE, quote=
```

Then, to perform pathway analysis online go to the Reactome website (<https://reactome.org/PathwayBrowser/#>

Q: What pathway has the most significant “Entities p-value”? Do the most significant pathways listed match your previous KEGG results? What factors could cause differences between the two methods?

The most significant entities p-value was the mitotic cell cycle. The most significant pathways mostly match the previous KEGG results. The differences can be caused due to variations in database structure, pathway coverage, and gene annotations. KEGG focuses on metabolic and signaling pathways, while Reactome provides more detailed molecular interactions and hierarchical biological processes.

GO online (optional)

Q: What pathway has the most significant “Entities p-value”? Do the most significant pathways listed match your previous KEGG results? What factors could cause differences between the two methods?

The pathway with the most significant p-value is the cell cycle. This matches the previous KEGG results but it also includes the trachea formation pathway. The differences can arise due to differences in databases and gene annotations.