

# Class 14: RNA-Seq analysis mini-project

Carly Chang (A16843962)

## Table of contents

|                                   |    |
|-----------------------------------|----|
| Background . . . . .              | 1  |
| Data Import . . . . .             | 2  |
| Inspect and tidy data . . . . .   | 2  |
| Setup for DESeq . . . . .         | 4  |
| Run DESeq . . . . .               | 4  |
| Volcano plot of results . . . . . | 5  |
| Gene annotation . . . . .         | 6  |
| Pathway Analysis . . . . .        | 8  |
| Gene Ontology Analysis . . . . .  | 13 |

## 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)
head(counts)
```

|                 | length    | SRR493366 | SRR493367 | SRR493368 | SRR493369 | SRR493370 |
|-----------------|-----------|-----------|-----------|-----------|-----------|-----------|
| ENSG00000186092 | 918       | 0         | 0         | 0         | 0         | 0         |
| ENSG00000279928 | 718       | 0         | 0         | 0         | 0         | 0         |
| ENSG00000279457 | 1982      | 23        | 28        | 29        | 29        | 28        |
| ENSG00000278566 | 939       | 0         | 0         | 0         | 0         | 0         |
| ENSG00000273547 | 939       | 0         | 0         | 0         | 0         | 0         |
| ENSG00000187634 | 3214      | 124       | 123       | 205       | 207       | 212       |
|                 | SRR493371 |           |           |           |           |           |
| ENSG00000186092 | 0         |           |           |           |           |           |
| ENSG00000279928 | 0         |           |           |           |           |           |
| ENSG00000279457 | 46        |           |           |           |           |           |
| ENSG00000278566 | 0         |           |           |           |           |           |
| ENSG00000273547 | 0         |           |           |           |           |           |
| ENSG00000187634 | 258       |           |           |           |           |           |

```
colData = read.csv("GSE37704_metadata.csv")
head(colData)
```

|   | 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      |

## Inspect and tidy data

Does the counts columns match the colData row?

```
colData$id
```

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

```
colnames(counts)
```

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

No, the two do not match. The fix here looks to be removing the first “length” column from counts:

```
countData <- counts[,-1] #gives everything except first column  
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) == colData$id
```

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

Q1. How many genes in total?

```
nrow(countData)
```

```
[1] 19808
```

Q2. Filter to remove zero count genes (rows where there are zero counts in all columns). How many genes are left?

```
new.countData <- countData[rowSums(countData) > 0,] #row indices where the sum is greater than 0  
head(new.countData)
```

|                 | 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.countData)
```

```
[1] 15975
```

## Setup for DESeq

```
library(DESeq2)
```

Setup input object for DESeq

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

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

## Run DESeq

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

estimating size factors

estimating dispersions

gene-wise dispersion estimates

mean-dispersion relationship

final dispersion estimates

fitting model and testing

```
res <- results(dds)
```

```
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>   |
| ENSG00000279457 | 29.9136     | 0.1792571      | 0.3248216 | 0.551863   | 5.81042e-01 |
| ENSG00000187634 | 183.2296    | 0.4264571      | 0.1402658 | 3.040350   | 2.36304e-03 |
| ENSG00000188976 | 1651.1881   | -0.6927205     | 0.0548465 | -12.630158 | 1.43989e-36 |
| ENSG00000187961 | 209.6379    | 0.7297556      | 0.1318599 | 5.534326   | 3.12428e-08 |
| ENSG00000187583 | 47.2551     | 0.0405765      | 0.2718928 | 0.149237   | 8.81366e-01 |
| ENSG00000187642 | 11.9798     | 0.5428105      | 0.5215599 | 1.040744   | 2.97994e-01 |
|                 | padj        |                |           |            |             |
|                 | <numeric>   |                |           |            |             |
| ENSG00000279457 | 6.86555e-01 |                |           |            |             |
| ENSG00000187634 | 5.15718e-03 |                |           |            |             |
| ENSG00000188976 | 1.76549e-35 |                |           |            |             |
| ENSG00000187961 | 1.13413e-07 |                |           |            |             |
| ENSG00000187583 | 9.19031e-01 |                |           |            |             |
| ENSG00000187642 | 4.03379e-01 |                |           |            |             |

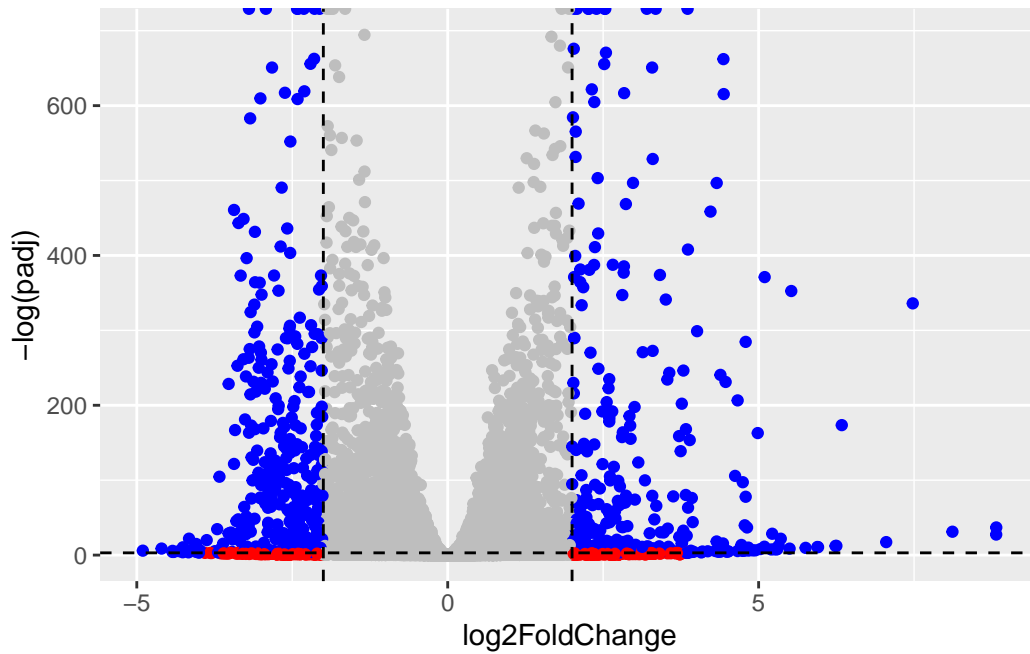
## Volcano plot of results

```
library(ggplot2)
```

```
mycols <- rep("grey", nrow(res))
mycols[abs(res$log2FoldChange) > 2] <- "blue"
mycols[abs(res$log2FoldChange) > 2 & res$padj > 0.05] <- "red"
```

```
ggplot(res) +
  aes(log2FoldChange, -log(padj)) +
  geom_point(col=mycols) +
  geom_vline(xintercept=-2,linetype=2) +
  geom_vline(xintercept=2,linetype=2) +
  geom_hline(yintercept=-log(0.05),linetype=2)
```

Warning: Removed 1237 rows containing missing values or values outside the scale range (`geom\_point()`).



## Gene annotation

Add gene SYMBOL and ENTREZID

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

```
columns(org.Hs.eg.db)
```

```
[1] "ACCNUM"      "ALIAS"       "ENSEMBL"     "ENSEMBLPROT" "ENSEMBLTRANS"
[6] "ENTREZID"    "ENZYME"      "EVIDENCE"     "EVIDENCEALL"  "GENENAME"
[11] "GENETYPE"    "GO"          "GOALL"        "IPI"           "MAP"
[16] "OMIM"        "ONTOLOGY"    "ONTOLOGYALL"  "PATH"          "PFAM"
[21] "PMID"        "PROSITE"     "REFSEQ"       "SYMBOL"        "UCSCCKG"
[26] "UNIPROT"
```

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

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

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

'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 8 columns

|                 | baseMean  | log2FoldChange | lfcSE     | stat       | pvalue      |
|-----------------|-----------|----------------|-----------|------------|-------------|
|                 | <numeric> | <numeric>      | <numeric> | <numeric>  | <numeric>   |
| ENSG00000279457 | 29.9136   | 0.1792571      | 0.3248216 | 0.551863   | 5.81042e-01 |
| ENSG00000187634 | 183.2296  | 0.4264571      | 0.1402658 | 3.040350   | 2.36304e-03 |
| ENSG00000188976 | 1651.1881 | -0.6927205     | 0.0548465 | -12.630158 | 1.43989e-36 |
| ENSG00000187961 | 209.6379  | 0.7297556      | 0.1318599 | 5.534326   | 3.12428e-08 |
| ENSG00000187583 | 47.2551   | 0.0405765      | 0.2718928 | 0.149237   | 8.81366e-01 |
| ENSG00000187642 | 11.9798   | 0.5428105      | 0.5215599 | 1.040744   | 2.97994e-01 |
|                 | padj      | symbol         | entrez    |            |             |

|                 | <numeric>   | <character> | <character> |
|-----------------|-------------|-------------|-------------|
| ENSG00000279457 | 6.86555e-01 | NA          | NA          |
| ENSG00000187634 | 5.15718e-03 | SAMD11      | 148398      |
| ENSG00000188976 | 1.76549e-35 | NOC2L       | 26155       |
| ENSG00000187961 | 1.13413e-07 | KLHL17      | 339451      |
| ENSG00000187583 | 9.19031e-01 | PLEKHN1     | 84069       |
| ENSG00000187642 | 4.03379e-01 | PERM1       | 84808       |

## Pathway Analysis

```
library(gage)
```

```
library(gageData)
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).
```

```
#####
```

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

Note that we used the mapIDs() function above to obtain Entrez gene IDs (stored in res\$entrez) and we have the fold change results from DESeq2 analysis (stored in res\$log2FoldChange).

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

| <NA>       | 148398     | 26155       | 339451     | 84069      | 84808      |
|------------|------------|-------------|------------|------------|------------|
| 0.17925708 | 0.42645712 | -0.69272046 | 0.72975561 | 0.04057653 | 0.54281049 |



Load up kegg dataset:

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

Run pathway analysis with KEGG:

```
keggres = gage(foldchanges, gsets=kegg.sets.hs)
head(keggres$less, 3)
```

|  | p.geomean    | stat.mean    |
|--|--------------|--------------|
| hsa04110 Cell cycle                            | 8.995727e-06 | -4.378644    |
| hsa03030 DNA replication                       | 9.424076e-05 | -3.951803    |
| hsa05130 Pathogenic Escherichia coli infection | 1.405864e-04 | -3.765330    |
|  | p.val        | q.val        |
| hsa04110 Cell cycle                            | 8.995727e-06 | 0.001889103  |
| hsa03030 DNA replication                       | 9.424076e-05 | 0.009841047  |
| hsa05130 Pathogenic Escherichia coli infection | 1.405864e-04 | 0.009841047  |
|  | set.size     | exp1         |
| hsa04110 Cell cycle                            | 121          | 8.995727e-06 |
| hsa03030 DNA replication                       | 36           | 9.424076e-05 |
| hsa05130 Pathogenic Escherichia coli infection | 53           | 1.405864e-04 |

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

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

Info: Working in directory /Users/carlychang/Downloads/UCSD/Winter 2025/BIMM 143/Class14

Info: Writing image file hsa04110.pathview.png

```
pathview(gene.data=foldchanges, pathway.id="hsa03030")
```

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

Info: Working in directory /Users/carlychang/Downloads/UCSD/Winter 2025/BIMM 143/Class14

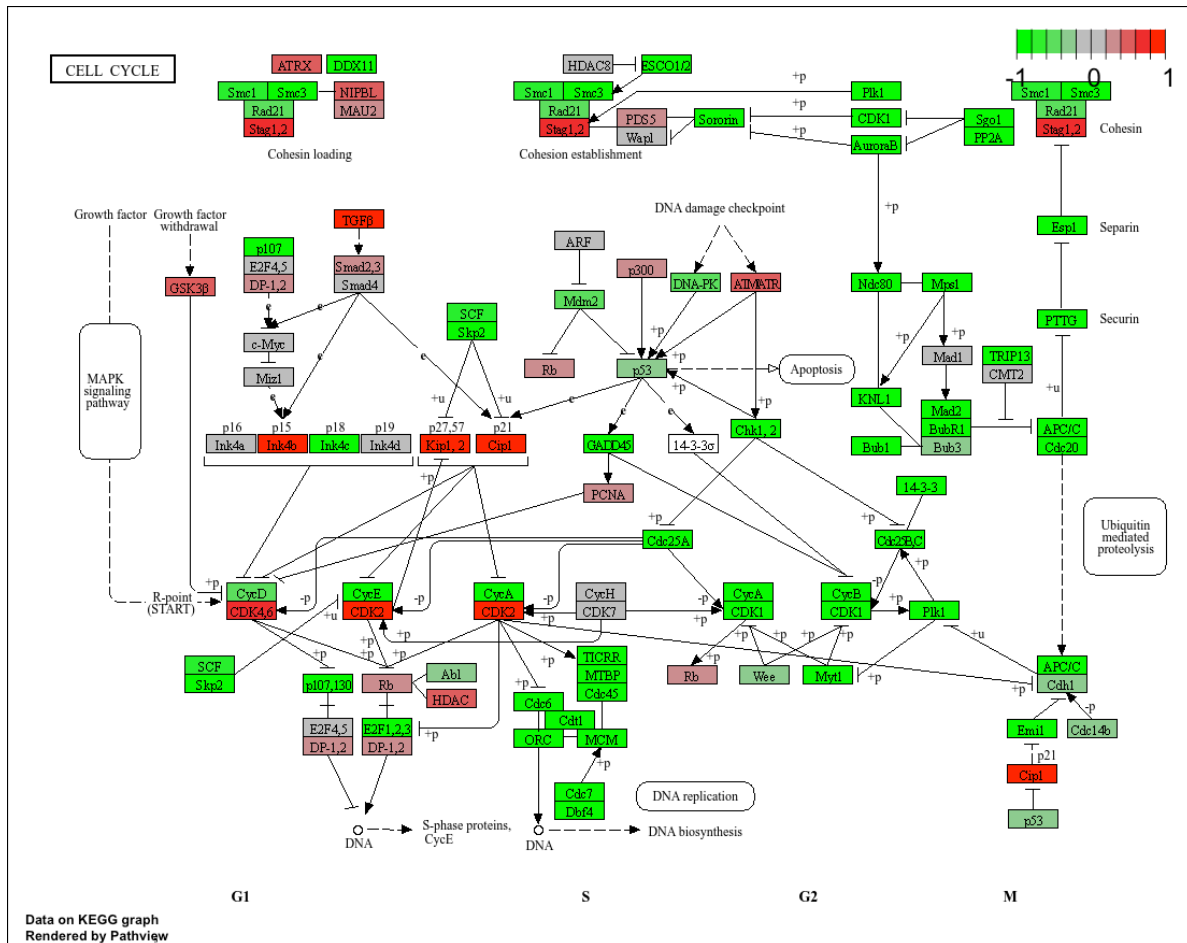
Info: Writing image file hsa03030.pathview.png

```
pathview(gene.data=foldchanges, pathway.id="hsa05130")
```

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

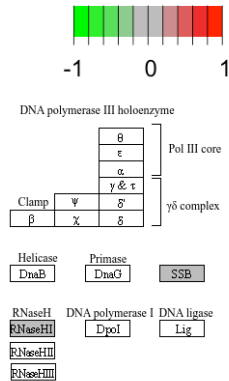
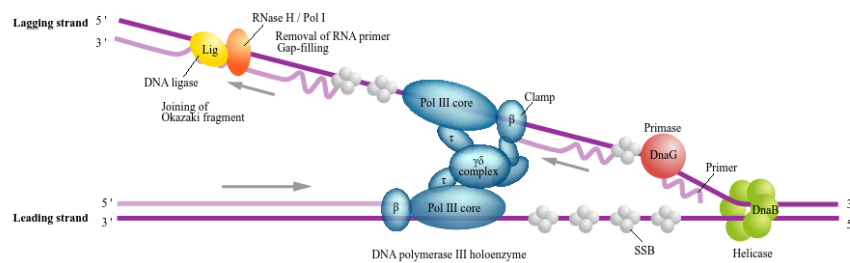
Info: Working in directory /Users/carlychang/Downloads/UCSD/Winter 2025/BIMM 143/Class14

Info: Writing image file hsa05130.pathview.png

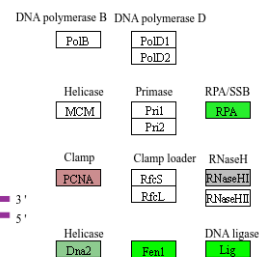
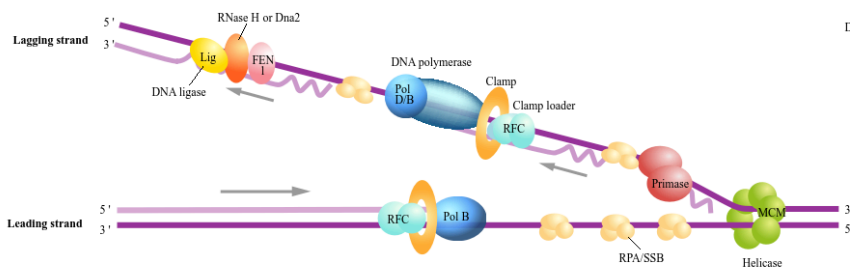


# DNA REPLICATION

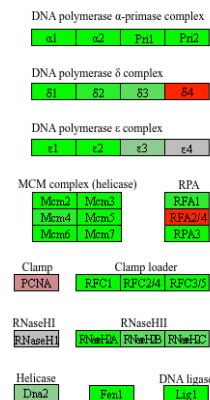
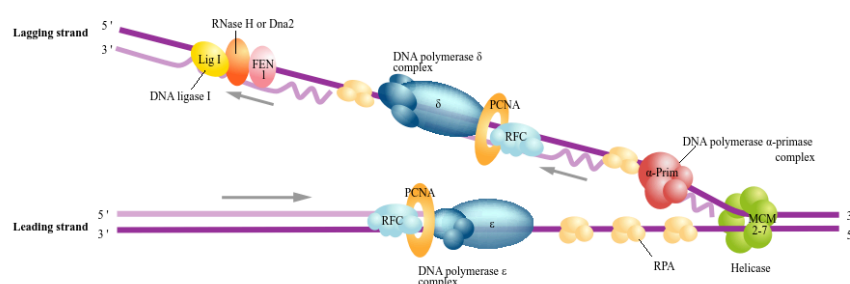
## Replication complex (Bacteria)



## Replication complex (Archaea)



## Replication complex (Eukaryotes)



Data on KEGG graph  
Rendered by Pathview



## Gene Ontology Analysis

Run pathway analysis with 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)

head(gobpres$less)
```

|  | p.geomean    | stat.mean | p.val        |
|--|--------------|-----------|--------------|
| G0:0048285 organelle fission             | 1.536227e-15 | -8.063910 | 1.536227e-15 |
| G0:0000280 nuclear division              | 4.286961e-15 | -7.939217 | 4.286961e-15 |
| G0:0007067 mitosis                       | 4.286961e-15 | -7.939217 | 4.286961e-15 |
| G0:0000087 M phase of mitotic cell cycle | 1.169934e-14 | -7.797496 | 1.169934e-14 |
| G0:0007059 chromosome segregation        | 2.028624e-11 | -6.878340 | 2.028624e-11 |
| G0:0000236 mitotic prometaphase          | 1.729553e-10 | -6.695966 | 1.729553e-10 |

|  | q.val        | set.size | expl         |
|--|--------------|----------|--------------|
| G0:0048285 organelle fission             | 5.841698e-12 | 376      | 1.536227e-15 |
| G0:0000280 nuclear division              | 5.841698e-12 | 352      | 4.286961e-15 |
| G0:0007067 mitosis                       | 5.841698e-12 | 352      | 4.286961e-15 |
| G0:0000087 M phase of mitotic cell cycle | 1.195672e-11 | 362      | 1.169934e-14 |
| G0:0007059 chromosome segregation        | 1.658603e-08 | 142      | 2.028624e-11 |
| G0:0000236 mitotic prometaphase          | 1.178402e-07 | 84       | 1.729553e-10 |