

# The Bio-Data Club @ Moffitt

Sponsored by the  
 R consortium

**WHO WE ARE:** Hosted by the Moffitt Biostatistics and Bioinformatics Shared Resource, the club is an informal group at Moffitt dedicated to promoting a fun, supportive, and hands-on environment for learning biological data science. Our meetings are intended for all levels of skill and computational backgrounds.

## PREVIOUS MEETINGS HAVE FEATURED...

Reproducible research with git and GitHub

Building interactive web applications with R/Shiny

Introduction to t-SNE clustering



**March 2019 meeting topic:**  
Automate your pathway enrichment analysis with R

**Date/time:** Friday, March 15th @ 2pm

**Location:** Moffitt Stabile Research Building (SRB), David Murphey Conference Room (1st floor)

**We meet once a month on Fridays. Bring a laptop! For a complete schedule and contact information, please visit the club GitHub page at <https://pstew.github.io/biodataclub>**



**MOFFITT**  
CANCER CENTER



# Automate your pathway enrichment analysis with R

2019-03-14

Paul Stewart, PhD

Bioinformatics Staff Scientist, BBSR

Sponsored by



**consortium**

## Why pathway analysis?

- Omics experiments can result in large lists of genes/proteins
  - Immediate relevance unclear
  - Relationship to one another unclear
- Genes can be grouped according to biological function, cellular compartment, chromosome, etc.
- Pathways help condense gene lists and prevent information overload.

**Big list of genes**

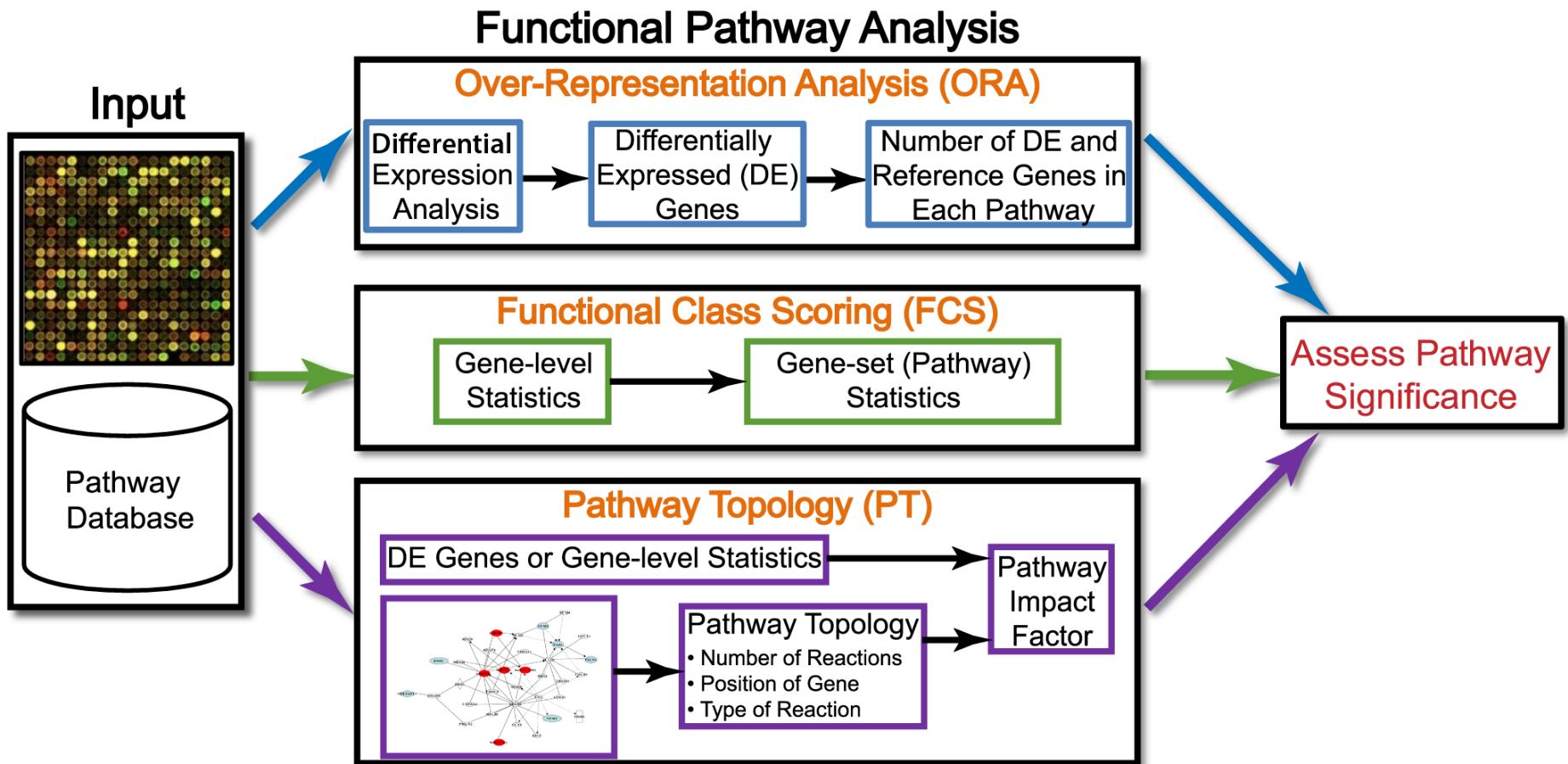
1  
2  
3  
4  
. . .  
. . .  
. . .  
. . .  
1,000



**Small list of pathways**

A  
B  
C  
D

Figure 1. Overview of existing pathway analysis methods using gene expression data as an example.



Khatri P, Sirota M, Butte AJ (2012) Ten Years of Pathway Analysis: Current Approaches and Outstanding Challenges. *PLOS Computational Biology* 8(2): e1002375. <https://doi.org/10.1371/journal.pcbi.1002375>

<http://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1002375>

## Over Representation Analysis

- In an experiment, you measured the expression of 30 genes.
- 15 differentially expressed genes were identified,
- 15 genes were associated with the GO pathway “DNA-templated transcription, elongation”.
- 12 genes are both differentially expressed and tagged with “translation initiation”.
  - Is this statistically significant?

Contingency table for gene expression data

	Differential Expression	NO Differential Expression	Total
IN Transcription Elongation	12	3	15
NOT IN Transcription Elongation	3	12	15
Total	15	15	30

# Over Representation Analysis Math (Fisher's Exact Test)

Contingency table for gene expression data

	Differential Expression	NO Differential Expression	Total
IN Transcription Elongation	12	3	15
NOT IN Transcription Elongation	3	12	15
Total	15	15	30

$$\text{"n choose k"} = \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

What is the probability for the table to the left?

$$P = \frac{(15 \text{ choose } 12) * (15 \text{ choose } 3)}{(30 \text{ choose } 15)}$$

$$= 1.33E-03$$

	C1	C2	C3	C4
R1	0   15 15   0 $p = 6.45E-09$	1   14 14   1 $p = 1.45E-06$	2   13 13   2 $p = 7.11E-05$	3   12 12   3 $p = 1.33E-03$
R2	4   11 11   4 $p = 1.20E-02$	5   10 10   5 $p = 5.81E-02$	6   9 9   6 $p = 1.61E-02$	7   8 8   7 $p = 2.67E-02$
R3	8   7 7   8 $p = 2.67E-02$	9   6 6   9 $p = 1.61E-02$	10   5 5   10 $p = 5.81E-02$	11   4 4   11 $p = 1.20E-02$
R4	12   3 3   12 $p = 1.33E-03$	13   2 2   13 $p = 7.11E-05$	14   1 1   14 $p = 1.45E-06$	15   0 0   15 $p = 6.45E-09$

Fisher's Exact Test amounts to summing the probability of observing our table of observed joint values **in addition to those more extreme than our table**, so

$$\text{P-value} = (R4, C1) + (R4, C2) + (R4, C3) + (R4, C4) = 0.0014$$

From this result, we claim that the probability of our observed data or that more extreme under the assumption that there is no association between expression and gene set membership is 0.0014.

## Input data

Choose an input file to upload. Either in BED format or a list of genes. For a quantitative set, add a comma and the level of membership of that gene. The membership level is a number between 0.0 and 1.0 to represent a weight for each gene, where the weight of 0.0 will completely discard the gene from the enrichment analysis and the weight of 1.0 is the maximum.

Try an example [BED file](#).

No file chosen

Or paste in a list of gene symbols optionally followed by a comma and levels of membership. Try two examples:  
[crisp set example](#), [fuzzy set example](#)

TALDO1  
ALDH3A1  
DDA1  
GSTA2  
PCYT1A  
GPC1  
GSTM3  
GSTM5  
PFN2  
PFN2

93 gene(s) entered

Enter a brief description for the list in case you want to share it. (Optional)

[Contribute](#)

Please acknowledge Enrichr in your publications by citing the following references:

Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, Clark NR, Ma'ayan A. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics*. 2013;128(14).

Kuleshov MV, Jones MR, Rouillard AD, Fernandez NF, Duan Q, Wang Z, Koplev S, Jenkins SL, Jagodnik KM, Lachmann A, McDermott MG, Monteiro CD, Gundersen GW, Ma'ayan A. Enrichr: a comprehensive gene set enrichment analysis web server 2016 update. *Nucleic Acids Research*. 2016; gkw377.

<http://amp.pharm.mssm.edu/Enrichr/>

Description Sample gene list (90 genes)



## KEGG 2016

Metabolism of xenobiotics by cytochrome P  
Chemical carcinogenesis\_Homo sapiens\_hsa  
Drug metabolism - cytochrome P450\_Homo  
Glutathione metabolism\_Homo sapiens\_hsa  
Metabolic pathways\_Homo sapiens\_hsa0110

## WikiPathways 2016

NRF2 pathway\_Homo sapiens\_WP2884  
Metapathway biotransformation\_Homo sapi  
Glutathione metabolism\_Homo sapiens\_WP  
Glutathione metabolism\_Mus musculus\_WP  
Aryl Hydrocarbon Receptor Pathway\_Homo

## ARCHS4 Kinases Coexp

RIPK4\_human\_kinase\_ARCHS4\_coexpression  
EPHB3\_human\_kinase\_ARCHS4\_coexpression  
SLK\_human\_kinase\_ARCHS4\_coexpression  
STK24\_human\_kinase\_ARCHS4\_coexpression  
PTK6\_human\_kinase\_ARCHS4\_coexpression

## Reactome 2016

Metabolism\_Homo sapiens\_R-HSA-1430728  
Biological oxidations\_Homo sapiens\_R-HSA-  
Phase II conjugation\_Homo sapiens\_R-HSA-  
Glutathione conjugation\_Homo sapiens\_R-H  
Glucuronidation\_Homo sapiens\_R-HSA-1565

## BioCarta 2016

Oxidative Stress Induced Gene Expression V  
Extrinsic Prothrombin Activation Pathway\_H  
Cystic fibrosis transmembrane conductance  
Y branching of actin filaments\_Homo sapien  
Mechanism of Gene Regulation by Peroxisor

## Humancyc 2016

glutathione-mediated detoxification\_Homo s  
putrescine degradation III\_Homo sapiens\_P  
allopregnanolone biosynthesis\_Homo sapien  
pentose phosphate pathway\_Homo sapiens\_  
superpathway of tryptophan utilization\_Hon

## NCI-Nature 2016

Validated transcriptional targets of TAp63 iso  
Posttranslational regulation of adherens jun  
Validated transcriptional targets of deltaNp6  
Direct p53 effectors\_Homo sapiens\_67c3b75  
S1P1 pathway\_Homo sapiens\_7327884f-619

## Panther 2016

5-Hydroxytryptamine degradation\_Homo sa  
Pentose phosphate pathway\_Homo sapiens  
p53 pathway\_Homo sapiens\_P00059  
p53 pathway feedback loops 2\_Homo sapien  
Heterotrimeric G-protein signaling pathway-

## BioPlex 2017

GPR18  
CDK15  
FCF1  
PI4KA  
FRMD1

Description Sample gene list (90 genes)



### GO Cellular Component 2017b

type III intermediate filament (GO:0045098)  
glial filament (GO:0097426)  
lamin filament (GO:0005638)  
neurofilament (GO:0005883)  
keratin filament (GO:0045095)

### GO Biological Process 2017b

glutathione biosynthetic process (GO:00067)  
epidermal cell differentiation (GO:0009913)  
glutathione derivative biosynthetic process (GO:00068)  
trypanothione biosynthetic process (GO:00068)  
glutathione deglycation (GO:0036531)

### GO Molecular Function 2017b

glutathione transferase activity (GO:0004364)  
protein homodimerization activity (GO:00422)  
NAD(P)H:methyl-1,4-benzoquinone oxidoreductase activity (GO:0004365)  
p-benzoquinone reductase (NADPH) activity (GO:0004366)  
NADPH:quinone reductase activity (GO:0004367)

### MGI Mammalian Phenotype 2017

MP:0001236\_abnormal\_epidermis\_stratum\_corniferum  
MP:0001730\_embryonic\_growth\_arrest  
MP:0000764\_abnormal\_tongue\_epithelium  
MP:0003809\_abnormal\_hair\_shaft\_morphology  
MP:0000352\_decreased\_cell\_proliferation

### Human Phenotype Ontology

Palmoplantar keratoderma (HP:0000982)  
Anonychia (HP:0001798)  
Abnormal blistering of the skin (HP:0008066)  
Woolly hair (HP:0002224)  
Skin ulcer (HP:0200042)

### Jensen TISSUES

Gut  
BTO:0004410  
Mouth  
Blood\_platelet  
Bronchial\_epithelial\_cell

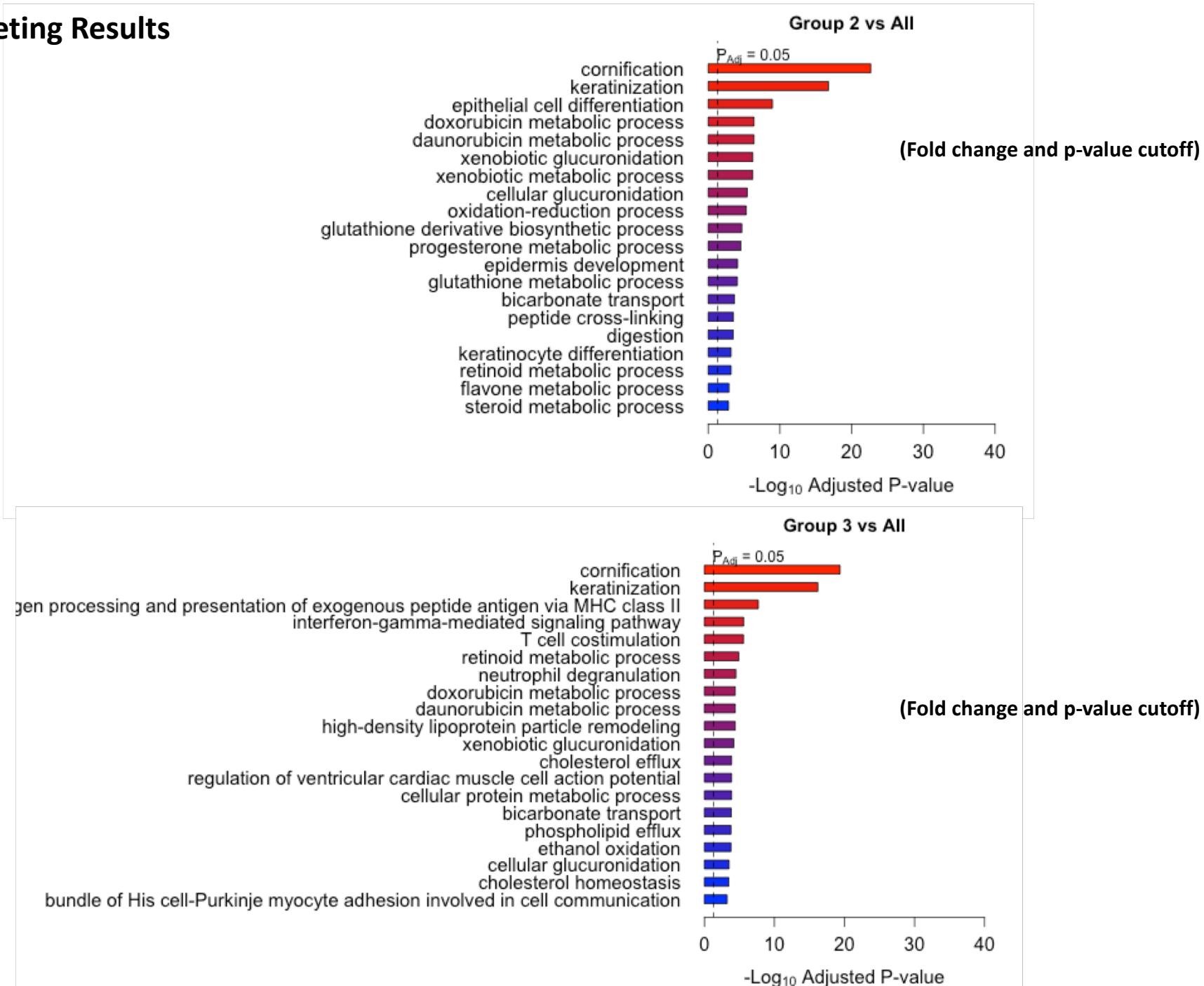
### Jensen COMPARTMENTS

Extracellular\_exosome  
Extracellular Vesicle  
Extracellular organelle  
Extracellular\_region\_part  
Extracellular\_region

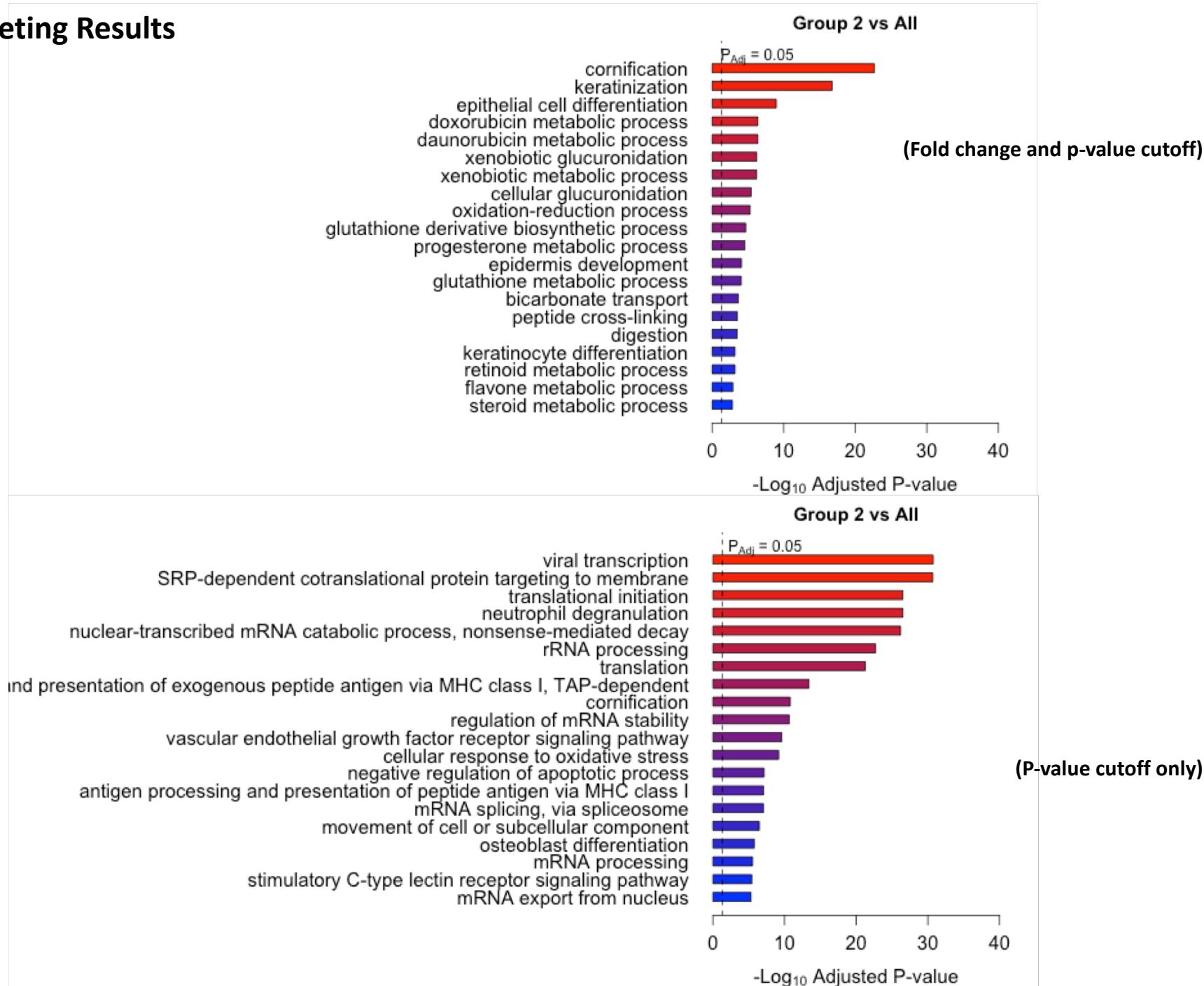
### Jensen DISEASES

Palmoplantar\_keratosis  
Anemia  
Pericholangitis  
Bilirubin\_metabolic\_disorder  
Fissured\_tongue

# Interpreting Results



# Interpreting Results



## Fold change and p-value cutoffs significantly alter microarray interpretations

Mark R Dalman<sup>1\*</sup>, Anthony Deeter<sup>2</sup>, Gayathri Nimishakavi<sup>2</sup>, Zhong-Hui Duan<sup>2</sup>

From Great Lakes Bioinformatics Conference 2011  
Athens, OH, USA. 2-4 May 2011

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2654802/>

BIOINFORMATICS

ORIGINAL PAPER

Vol. 25 no. 6 2009, pages 765–771  
doi:10.1093/bioinformatics/btp053

Gene expression

## Testing significance relative to a fold-change threshold is a TREAT

Davis J. McCarthy and Gordon K. Smyth\*

The Walter and Eliza Hall Institute of Medical Research, 1G Royal Parade, Parkville, Victoria 3050, Australia

Received on October 12, 2008; revised on January 21, 2009; accepted on January 22, 2009

Advance Access publication January 28, 2009

Associate Editor: Joaquin Dopazo

<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3305783/>