# Module 3, Lesson 17
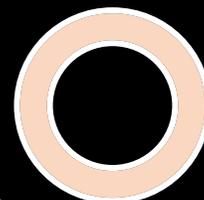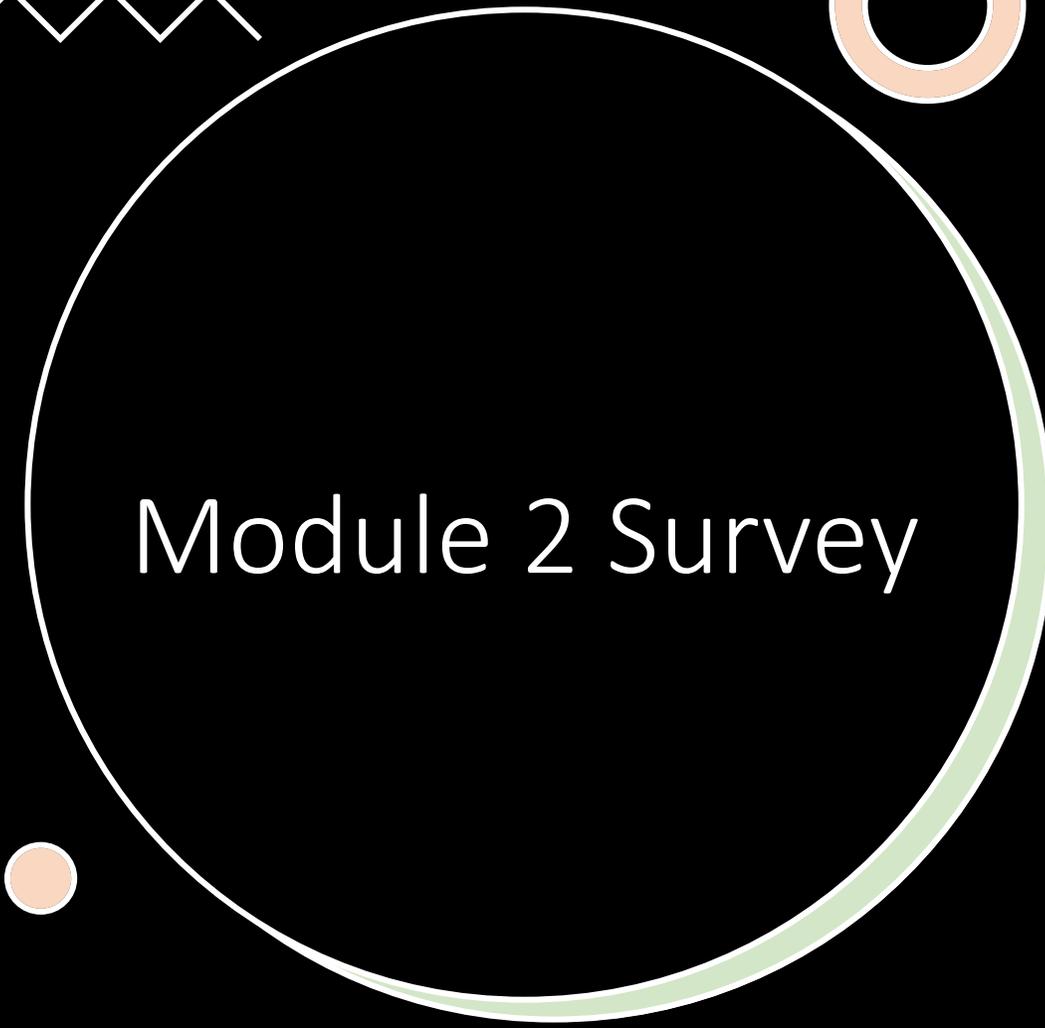# Introduction to Pathway Analysis

Alexandra Emmons, PhD
December 1, 2022

# Bioinformatics Training and Education Program

NIH NATIONAL CANCER INSTITUTE
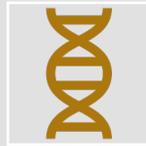Center for Cancer Research

# Module 2 Survey

- Let's take a second to take a brief Webex poll

- This will help us improve the class and clear up any misunderstandings by the final lesson.
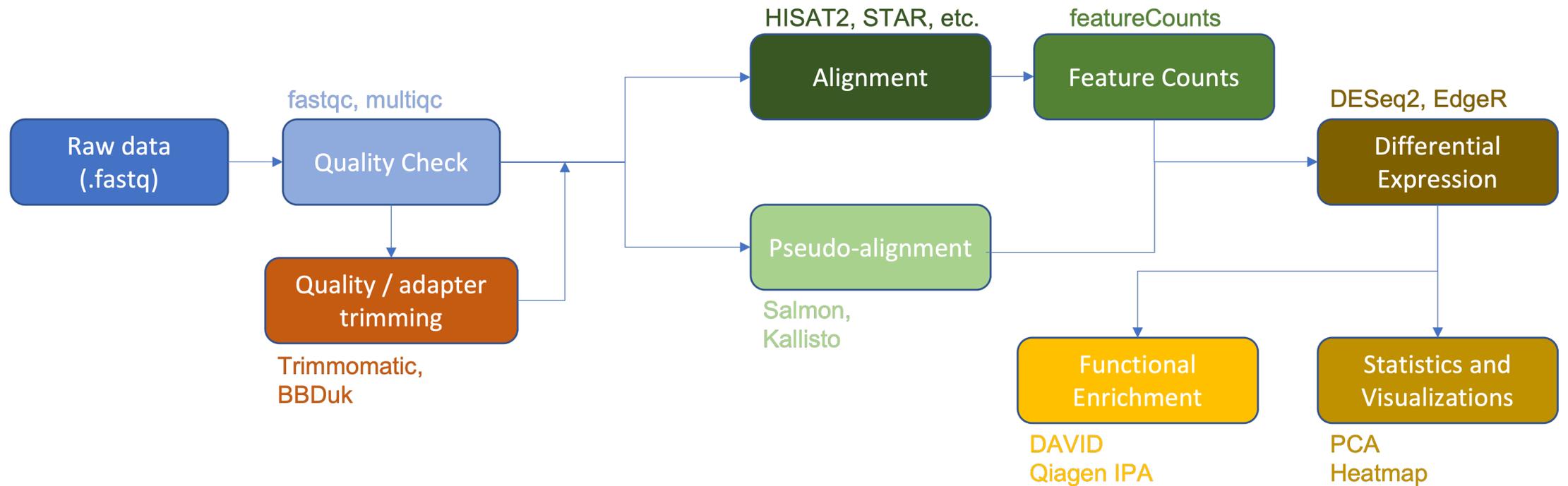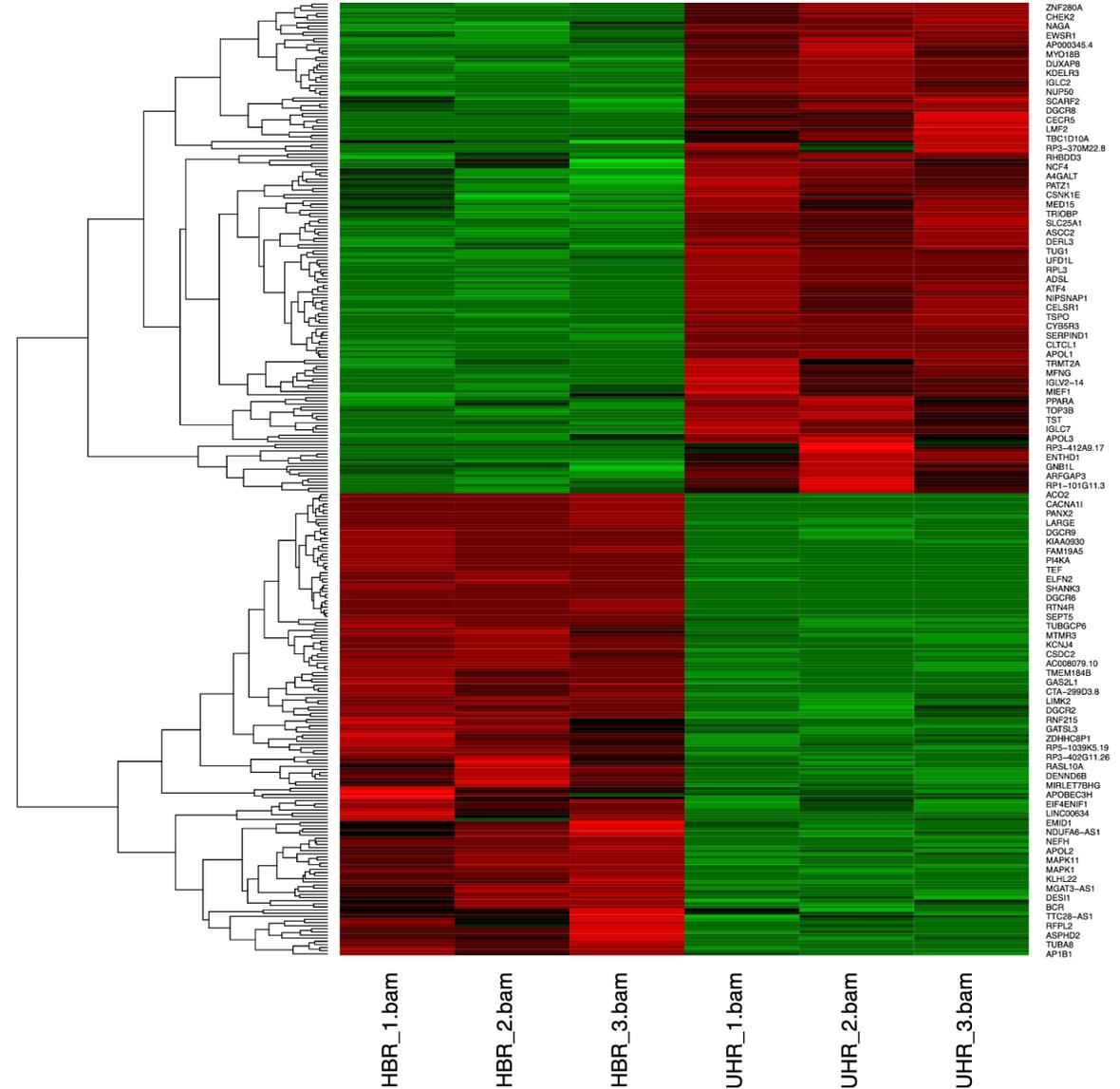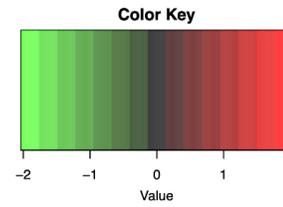
# Lesson 17 objectives

- Determine potential next steps following differential expression analysis.

- Tour geneontology.org and understand the three main ontologies.

- Learn about different methods and tools related to functional enrichment and pathway analysis.

- Get familiar with databases commonly used by popular functional enrichment tools.
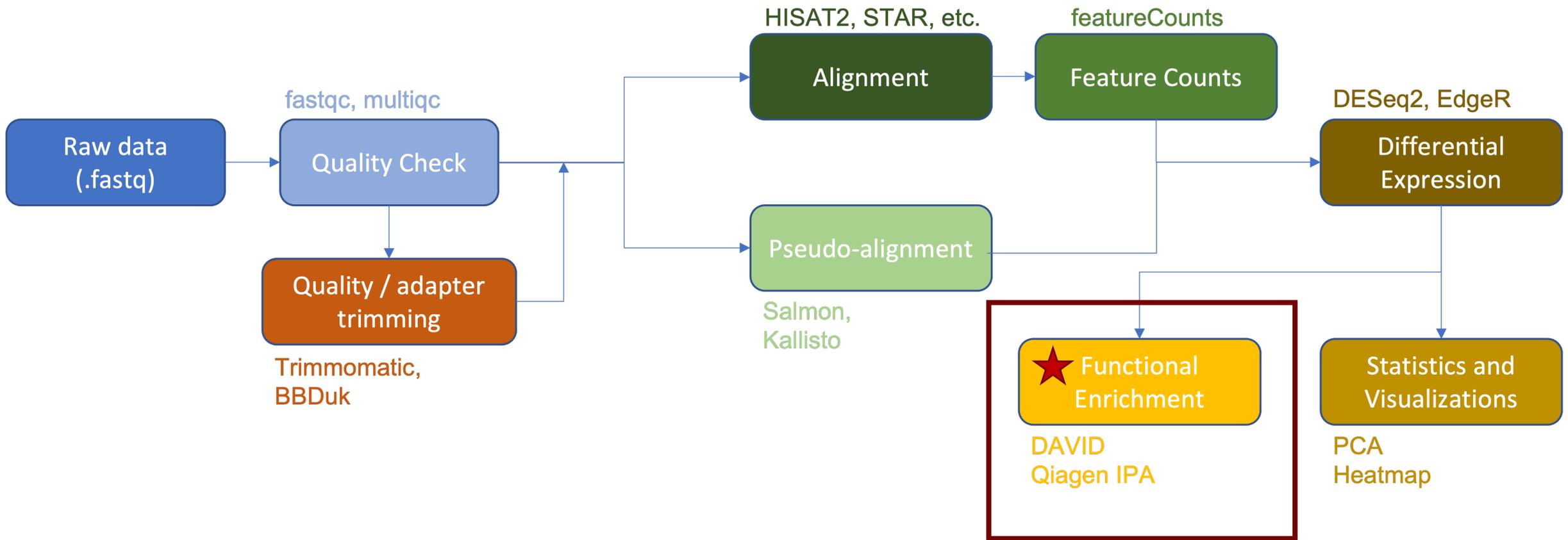
# Overview

Raw data (.fastq) → Quality Check

fastqc, multiqc

Quality Check → Quality / adapter trimming

Trimmomatic, BBDuk

HISAT2, STAR, etc.

Alignment

featureCounts

Feature Counts

Salmon, Kallisto

Pseudo-alignment

DESeq2, EdgeR

Differential Expression

Functional Enrichment

DAVID
Qiagen IPA

Statistics and Visualizations

PCA
Heatmap

# SO MANY GENES....

We generated a heatmap, so what's next?

# Why gene set / pathway analysis?

**1** Increase the statistical power in our analysis

**2** Ease interpretation

**3** Predict new roles for genes

**4** Better integrate data from different methods

# What is gene ontology?

- The Gene Ontology (GO) provides a framework and set of concepts for describing the functions of gene products from all organisms. --- https://www.ebi.ac.uk/ols/ontologies/go
  - Controlled vocabulary
  - Maintained by the Gene Ontology Consortium
    - Updated regularly
- Two parts:
  - the ontology (the GO terms and their hierarchical relationship)
  - the annotations (the annotated genes linked to various GO terms)

# What is gene ontology?

GO integrates information about gene product function in the context of three domains:

- Molecular function (MF) - "the molecular activities of individual gene products"
- Cellular component (CC) - "where the gene products are active"
- Biological process (BP) - "the pathways and larger processes to which that gene product's activity contributes"

QuickGO - https://www.ebi.ac.uk/QuickGO

# GO Statistics

## Ontology

| Property | Value |
| --- | --- |
| Valid terms | 43303 (Δ = -26) |
| Obsoleted terms | 4094 (Δ = 71) |
| Merged terms | 2442 (Δ = 4) |
| Biological process terms | 27993 |
| Molecular function terms | 11271 |
| Cellular component terms | 4039 |

## Annotations

| Property | Value |
| --- | --- |
| Number of annotations | 7,687,289 |
| Annotations for biological process | 2,872,350 |
| Annotations for molecular function | 2,432,692 |
| Annotations for cellular component | 2,382,247 |
| Annotations for evidence PHYLO | 3,993,931 |
| Annotations for evidence IEA | 1,573,469 |
| Annotations for evidence OTHER | 871,395 |
| Annotations for evidence EXP | 937,340 |
| Annotations for evidence ND | 252,104 |
| Annotations for evidence HTP | 59,050 |
| Number of annotated scientific publications | 172,927 |

## Gene products and species

| Property | Value |
| --- | --- |
| Annotated gene products | 1,503,630 |
| Annotated species | 5,257 |
| Annotated species with over 1,000 annotations | 185 |

Check out these tips for working with GO terms!

Current release 2022-11-03: 43,303 GO terms | 7,687,289 annotations
1,503,630 gene products | 5,257 species  (see statistics)

# THE GENE ONTOLOGY RESOURCE

The mission of the GO Consortium is to develop a comprehensive, **computational model of biological systems**, ranging from the molecular to the organism level, across the multiplicity of species in the tree of life.

The Gene Ontology (GO) knowledgebase is the world's largest source of information on the functions of genes. This knowledge is both human-readable and machine-readable, and is a foundation for computational analysis of large-scale molecular biology and genetics experiments in biomedical research.

Search GO term or Gene Product in AmiGO ...

Hint: add a space after completing a word to narrow the search.

● Any    ● Ontology    ● Gene Product

## GO Enrichment Analysis ❓

*Powered by PANTHER*

Your gene IDs here...

biological process

Homo sapiens        Examples    Launch ❯

*Hint: can use UniProt ID/AC, Gene Name, Gene Symbols, MOD IDs*

---

### 🔗 ONTOLOGY

The network of biological classes describing the current best representation of the "universe" of biology: the molecular functions, cellular locations, and processes gene products may carry out.

ⓘ GO Ontology Overview
⧉ Browse in AmiGO
⬇ Download

### ✎ ANNOTATION

NEDD4 -------- (evidence) -------- Ubiquitin-protein ligase activity GO:0004842

Statements, based on specific, traceable scientific evidence, asserting that a specific gene product is a real exemplar of a particular GO class.

ⓘ GO Annotations Overview
⧉ Browse in AmiGO
⬇ Download

### ⚙ GO-CAM

GO Causal Activity Model (GO-CAM) provides a structured framework to link standard GO annotations into a more complete model of a biological system.

ⓘ GO-CAM Overview
⧉ Browse GO-CAMs
⬇ Download

### ⚙ TOOLS & GUIDES

Tools to curate, browse, search, visualize and download both the ontology and annotations. Includes bioinformatic guides (Notebooks) and simple API access to integrate the GO into your research.

ⓘ GO Tools Overview
⧉ GO APIs Guide
⌥ GO GitHub

# Other databases

- Kyoto Encyclopedia of Genes and Genomes (KEGG)
  - Curated database
  - biological pathways
  - Molecular interaction networks
  - Very nice pathway maps
  - Restricted licenses
- Pathway Commons
  - a meta-database of pathways from other pathway databases
- PANTHER
  - Database of signaling pathways
- WikiPathways
  - community driven meta-database of pathways

# Other databases

- NDEx
  - an open-source framework where scientists and organizations can store, share, manipulate, and publish biological network knowledge

- HumanCyc
  - an encyclopedic reference on human metabolic pathways, the human genome, and human metabolites.

Check out Pathguide to get an idea of how many databases are available.

# 3 general approaches to pathway analysis:



Image from Khatri et al. 2012

- Over-representation analysis
- Functional class scoring
- Pathway topology

# Over-representation Analysis

- Statistically evaluates the fraction of genes in a particular pathway found among the set of genes showing changes in expression --- (Khatri et al. 2012)

- Strategy
  - Provide an input list of gene IDs (uses a threshold)
  - Input genes for each pathway are counted
  - The same counting step is applied to a background set of genes
  - Pathways are tested for over or under representation using tests based on hypergeometric, chi-square, or binomial distribution

# DAVID



- **D**atabase for **A**nnotation, **V**isualization and **I**ntegrated **D**iscovery

- Very popular

- Easy to use and produces a lot of output

- Uses a variety of databases (NCBI, Uniprot, Ensembl, Gene Ontology, KEGG, Reactome, etc.).

# Functional Class Scoring

- "The hypothesis of functional class scoring (FCS) is that although large changes in individual genes can have significant effects on pathways, weaker but coordinated changes in sets of functionally related genes (i.e., pathways) can also have significant effects." --- Khatri et al. 2012

- Strategy:
  - Compute a gene level statistic (differential expression)
  - Create a pathway level statistic by aggregating gene level stats
  - Determine statistical significance from pathway stat
    - Competitive vs self-contained methods

# Functional Class Scoring: GSEA



Image from https://diytranscriptomics.com/project/lecture-10



Fig 1: Enrichment plot: P53_DOWN_KANNAN
Profile of the Running ES Score & Positions of GeneSet Members on the Rank Ordered List

# Pathway Topology

- PT methods generally use network based modeling
- They consider the information ignored in ORA and FCS methods: gene product interactions, positions of genes, and types of genes
- Examples iPathwayGuide, Pathway-Express, SPIA, NetGSA

# Some possible tools

Gene set analysis tools

| Tool | Author | Year | Citations[1] | Availability | Gene sets | Methods[2] |
|---|---|---|---|---|---|---|
| WEBGESTALT | Zhang *et al.* [73] | 2005 | 1423 | Web server | GO, KEGG, +20 more | ORA, GSEA |
| GOSTATS | Falcon and Gentleman [74] | 2007 | 1437 | R package | GO | ORA |
| G:PROFILER | Reimand *et al.* [75] | 2007 | 534 | Web server | GO, KEGG, +7 more | ORA |
| GENETRAIL | Backes *et al.* [76] | 2007 | 360 | Web server | GO, KEGG, +28 more | ORA, GSEA |
| DAVID | Huang *et al.* [8] | 2009 | 19 569 | Web server | GO, KEGG, +38 more | ORA |
| GORILLA | Eden *et al.* [77] | 2009 | 1881 | Web server | GO | ORA |
| TOPPGENE | Chen *et al.* [78] | 2009 | 1200 | Web server | GO, KEGG, +45 more | ORA |
| CLUSTER-PROFILER | Yu *et al.* [10] | 2012 | 1305 | R package | GO, KEGG, +8 more | ORA, GSEA |
| PANTHER | Mi *et al.* [79] | 2013 | 1405 | Web server | GO, +2 more | ORA, GSEA |
| ENRICHR | Chen *et al.* [9] | 2013 | 1246 | Web server | GO, KEGG, +33 more | ORA |

[1]Google Scholar, July 2019.

[2]Detailed summary of implemented methods in Supplementary Methods S1.2.

# A note about tools and databases

- Not all databases survive
  - Check to see when information was last updated
- Tools also frequently fall to the wayside
  - Check to see if the tool is maintained
  - If the tool is not readily updated, it's likely it is using outdated versions of databases.

# Importance of gene IDs

- Different tools use different gene ID annotations
- May need to convert between annotations using available programs
- Annotations can be genome build specific
- See linked tutorial in the course docs
- Some annotation programs/databases:
  - g:convert
  - BioMart
  - AnnotationHub

# R packages

# Pathview

# Pathway enrichment analysis and visualization of omics data using g:Profiler, GSEA, Cytoscape and EnrichmentMap

Jüri Reimand[1,2,8], Ruth Isserlin[3,8], Veronique Voisin[3], Mike Kucera[3], Christian Tannus-Lopes[3], Asha Rostamianfar[3], Lina Wadi[1], Mona Meyer[1], Jeff Wong[3], Changjiang Xu[3], Daniele Merico[4,5] and Gary D. Bader[3,6,7*]

Pathway enrichment analysis helps researchers gain mechanistic insight into gene lists generated from genome-scale (omics) experiments. This method identifies biological pathways that are enriched in a gene list more than would be expected by chance. We explain the procedures of pathway enrichment analysis and present a practical step-by-step guide to help interpret gene lists resulting from RNA-seq and genome-sequencing experiments. The protocol comprises three major steps: definition of a gene list from omics data, determination of statistically enriched pathways, and visualization and interpretation of the results. We describe how to use this protocol with published examples of differentially expressed genes and mutated cancer genes; however, the principles can be applied to diverse types of omics data. The protocol describes innovative visualization techniques, provides comprehensive background and troubleshooting guidelines, and uses freely available and frequently updated software, including g:Profiler, Gene Set Enrichment Analysis (GSEA), Cytoscape and EnrichmentMap. The complete protocol can be performed in ~4.5 h and is designed for use by biologists with no prior bioinformatics training.