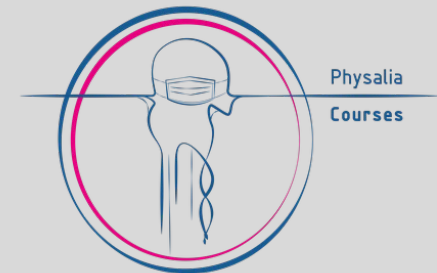# Gene ontology, gene set over-representation analyses and public databases
## Peeking into functional roles of gene sets

Epigenomics Data Analysis

Jacques Serizay

Physalia 2025

# Why do we need GO analyses?

After RNA-seq (or ATAC, ChIP, …), we generally end up with a list of genes (or genomic features with the nearest gene).

➔ **We want to investigate the biological relevance of this set of genes.**

An ontology is a **formal representation** of a **body of knowledge** **within a given domain**

**An ontology term** primarily consists of:

- A definition of a concept

- A representation of this concept

- A formal naming of this concept

# GO over-representation analyses: what

**An ontology term** primarily consists of:

- A definition of a concept

- A representation of this concept

- A formal naming of this concept

```
=== Example term ===
:id:            GO:0000016
:name:          lactase activity
:ontology:      molecular_function
:def:           "Catalysis of the reaction: lactose + H2O=D-
                glucose + D-galactose." [EC:3.2.1.108]
:synonym:        "lactase-phlorizin hydrolase activity"
                BROAD [EC:3.2.1.108]
:synonym:       "lactose galactohydrolase activity" EXACT
                [EC:3.2.1.108]
:xref:          EC:3.2.1.108
:xref:          MetaCyc:LACTASE-RXN
:xref:          Reactome:20536
:is_a:          GO:0004553 ! hydrolase activity,
                hydrolyzing O-glycosyl compounds
```
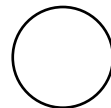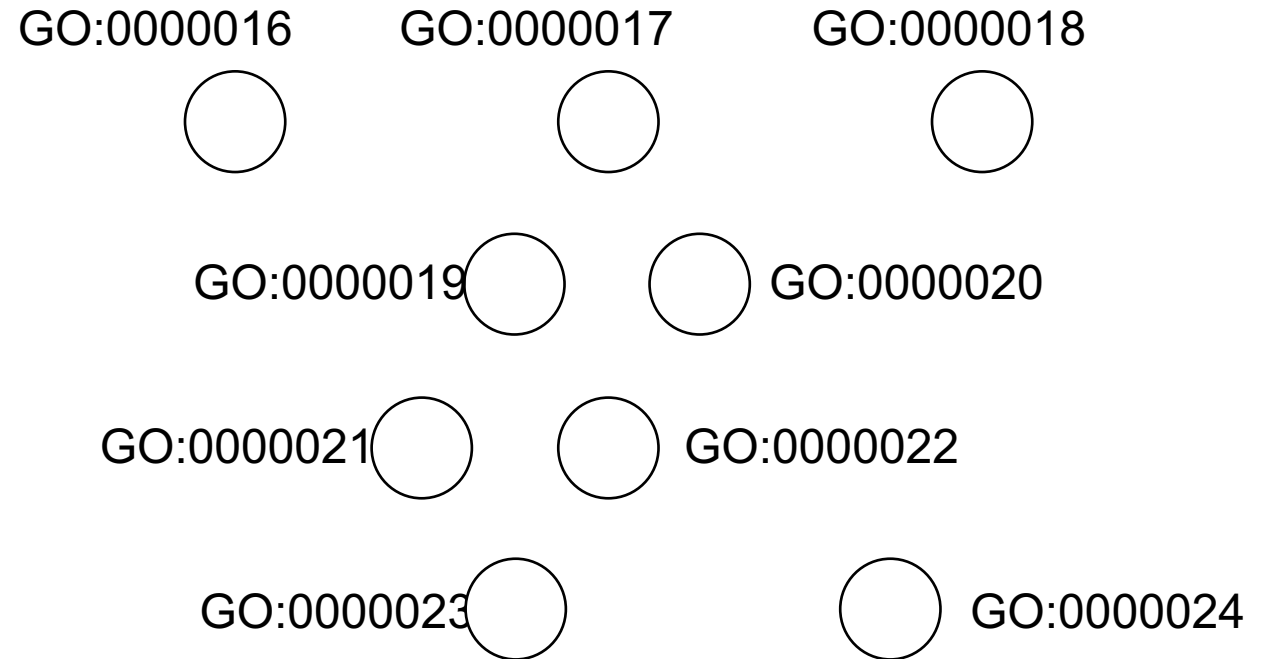
# GO over-representation analyses: what

**An ontology term** primarily consists of:

- A definition of a concept

- A representation of this concept

- A formal naming of this concept

```
=== Example term ===
:id:            GO:0000016
:name:          lactase activity
:ontology:      molecular_function
:def:           "Catalysis of the reaction: lactose + H2O=D-
                glucose + D-galactose." [EC:3.2.1.108]
:synonym:        "lactase-phlorizin hydrolase activity"
                BROAD [EC:3.2.1.108]
:synonym:       "lactose galactohydrolase activity" EXACT
                [EC:3.2.1.108]
:xref:          EC:3.2.1.108
:xref:          MetaCyc:LACTASE-RXN
:xref:          Reactome:20536
:is_a:          GO:0004553 ! hydrolase activity,
                hydrolyzing O-glycosyl compounds
```

**GO:0000016**

**An ontology term** can be further enriched with additional information:

- Elements can be annotated to individual terms

**B**
**D** GO:0000022

## Ontology terms



GO:0000016

GO:0000017

GO:0000018

GO:0000019

GO:0000020

GO:0000021

GO:0000022

GO:0000023

GO:0000024

**Gene Ontology terms**

**+ GO annotations**

**Gene Ontology terms**

**+ GO annotations**

**+ hierarchy**

Different **Ontology terms** can contain the same sets of annotations
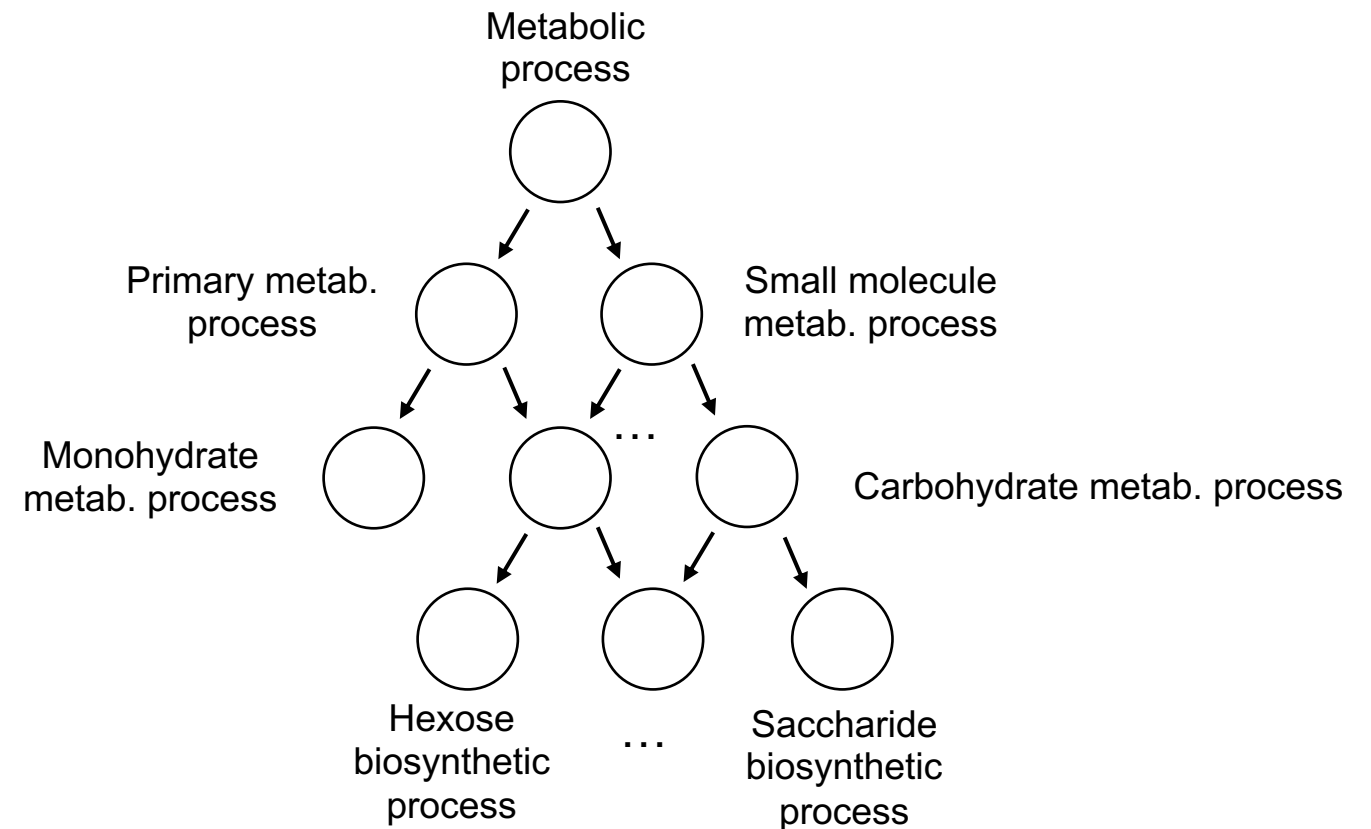
In our case, the Gene Ontology (GO) describes the **current state** of <u>knowledge of the three main biological domains</u>

In our case, the Gene Ontology (GO) describes the **current state** of <u>knowledge in biology</u>
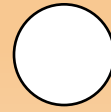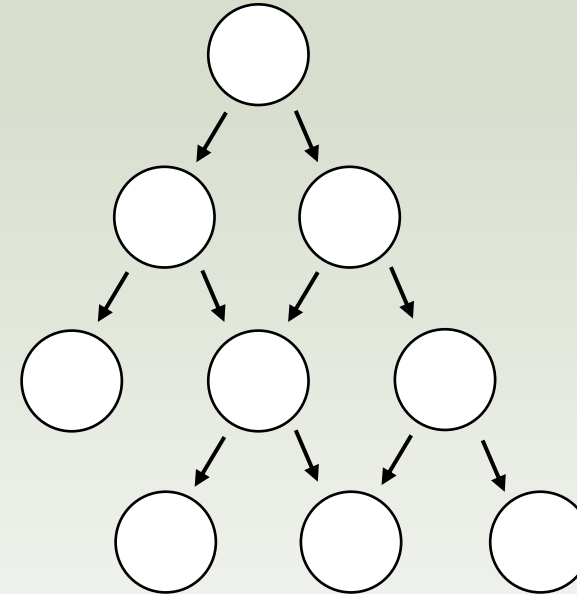
Gene Ontology (GO) is divided in three domains

- Biological Processes (BP)

- Cellular Components (CC)

- Molecular Functions (MF)

The Gene Ontology (GO) is a dynamic, frequently updated database

**IMPORTANT:**

A <u>GO term</u> (e.g. GO:0000017) is different from its

<u>annotations</u> (i.e. the association of some genes

to this term)

GO:0000017

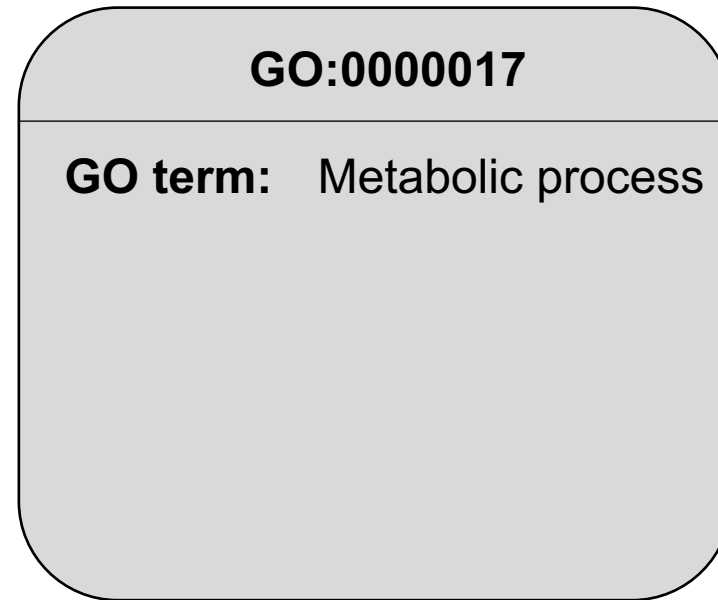| **GO:0000017** |
|:---|
| **GO term:**   Metabolic process |

## GO over-representation analyses: what

**IMPORTANT:**

A GO term (e.g. GO:0000017) is different from its

annotations (i.e. the association of some genes

to this term)

GO:0000017

A B
C D

**GO:0000017**

**GO term:**   Metabolic process

**GO annotations:**   Gene A
Gene B
Gene C
Gene D

I apologize — let me provide the clean transcription.

**IMPORTANT:**

**GO consortium** organizes **GO terms** and their hierarchy

**External providers** manage **GO term annotations**

**IMPORTANT:**

GO consortium organizes GO terms and their hierarchy

External providers manage GO term annotations

- o Mouse annotations are provided by MGI (Mouse Genome Informatics)

- o C. elegans annotations are provided by Wormbase

- o Yeast annotations are provided by SGD (Saccharomyces Genome Database)

# Official GO database

GO Consortium is the provider of official Gene Ontology.

# Downloading official GO database

- Versioned database

- Easy access to entire database

- OBO format

# Downloading official GO database

OBO format

## wget http://purl.obolibrary.org/obo/go.obo

```
 1   > head -n 100 go.obo
 2   format-version: 1.2
 3   data-version: releases/2020-12-08
 4   ontology: go
 5
 6   [Term]
 7   id: GO:0000001
 8   name: mitochondrion inheritance
 9   namespace: biological_process
10   def: "The distribution of mitochondria, including the mitochondrial genome, into daughter cells after mitosis or meiosis, mediated by interactions
       between mitochondria and the cytoskeleton." [GOC:mcc, PMID:10873824, PMID:11389764]
11   synonym: "mitochondrial inheritance" EXACT []
12   is_a: GO:0048308 ! organelle inheritance
13   is_a: GO:0048311 ! mitochondrion distribution
14
15   [Term]
16   id: GO:0000002
17   name: mitochondrial genome maintenance
18   namespace: biological_process
19   def: "The maintenance of the structure and integrity of the mitochondrial genome; includes replication and segregation of the mitochondrial chromosome."
       [GOC:ai, GOC:vw]
20   is_a: GO:0007005 ! mitochondrion organization
```

GO **annotations:**

- Also versioned

- All most recent annotations for individual species available at:

http://current.geneontology.org/annotations/

# Downloading GO annotations

```
Parent                             goa_dog_rna.gpad.gz              goa_chicken_complex.gpad.gz     mgi.gaf.gz
..                                 goa_dog_rna.gpi.gz               goa_chicken_complex.gpi.gz      mgi.gpad.gz
                                   goa_human.gaf.gz                 goa_chicken_isoform.gaf.gz      mgi.gpi.gz
                                   goa_human.gpad.gz                goa_chicken_isoform.gpad.gz     pombase.gaf.gz
                                   goa_human.gpi.gz                 goa_chicken_isoform.gpi.gz      pombase.gpad.gz
aspgd.gaf.gz                       goa_human_complex.gaf.gz         goa_chicken_rna.gaf.gz          pombase.gpi.gz
aspgd.gpad.gz                      goa_human_complex.gpad.gz        goa_chicken_rna.gpad.gz         pseudocap.gaf.gz
aspgd.gpi.gz                       goa_human_complex.gpi.gz         goa_chicken_rna.gpi.gz          pseudocap.gpad.gz
cgd.gaf.gz                         goa_human_isoform.gaf.gz         goa_cow.gaf.gz                  pseudocap.gpi.gz
cgd.gpad.gz                        goa_human_isoform.gpad.gz        goa_cow.gpad.gz                 reactome.gaf.gz
cgd.gpi.gz                         goa_human_isoform.gpi.gz         goa_cow.gpi.gz                  reactome.gpad.gz
dictybase.gaf.gz                   goa_human_rna.gaf.gz             goa_cow_complex.gaf.gz          reactome.gpi.gz
dictybase.gpad.gz                  goa_human_rna.gpad.gz            goa_cow_complex.gpad.gz         rgd.gaf.gz
dictybase.gpi.gz                   goa_human_rna.gpi.gz             goa_cow_complex.gpi.gz          rgd.gpad.gz
ecocyc.gaf.gz                      goa_pig.gaf.gz                   goa_cow_isoform.gaf.gz          rgd.gpi.gz
ecocyc.gpad.gz                     goa_pig.gpad.gz                  goa_cow_isoform.gpad.gz         sgd.gaf.gz
ecocyc.gpi.gz                      goa_pig.gpi.gz                   goa_cow_isoform.gpi.gz          sgd.gpad.gz
fb.gaf.gz                          goa_pig_complex.gaf.gz           goa_cow_rna.gaf.gz              sgd.gpi.gz
fb.gpad.gz                         goa_pig_complex.gpad.gz          goa_cow_rna.gpad.gz             sgn.gaf.gz
fb.gpi.gz                          goa_pig_complex.gpi.gz           goa_cow_rna.gpi.gz              sgn.gpad.gz
genedb_lmajor.gaf.gz               goa_pig_isoform.gaf.gz           goa_dog.gaf.gz                  sgn.gpi.gz
genedb_lmajor.gpad.gz              goa_pig_isoform.gpad.gz          goa_dog.gpad.gz                 tair.gaf.gz
genedb_lmajor.gpi.gz               goa_pig_isoform.gpi.gz           goa_dog.gpi.gz                  tair.gpad.gz
genedb_tbrucei.gaf.gz              goa_pig_rna.gaf.gz               goa_dog_complex.gaf.gz          tair.gpi.gz
genedb_tbrucei.gpad.gz             goa_pig_rna.gpad.gz              goa_dog_complex.gpad.gz         wb.gaf.gz
genedb_tbrucei.gpi.gz              goa_pig_rna.gpi.gz               goa_dog_complex.gpi.gz          wb.gpad.gz
goa_chicken.gaf.gz                 goa_uniprot_all.gaf.gz           goa_dog_isoform.gaf.gz          wb.gpi.gz
goa_chicken.gpad.gz                goa_uniprot_all_noiea.gaf.gz     goa_dog_isoform.gpad.gz         zfin.gaf.gz
goa_chicken.gpi.gz                 goa_uniprot_all_noiea.gpad.gz    goa_dog_isoform.gpi.gz          zfin.gpad.gz
goa_chicken_complex.gaf.gz         goa_uniprot_all_noiea.gpi.gz     goa_dog_rna.gaf.gz              zfin.gpi.gz
```

GAF format:


It's in the name:  **G**O **Annotation F**ormat

# Downloading GO annotations

GAF format:

```
!gaf-version: 2.1
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-09
!
!Header from source association file:
!===============================
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-08
!
!Header from sgd source association file:
!===============================
!Date: 20201207
!From: Saccharomyces Genome Database (SGD)
!URL: https://www.yeastgenome.org/
!Contact Email: sgd-helpdesk@lists.stanford.edu
!Funding: NHGRI at US NIH, grant number U41-HG001315
!
!===============================
!
!Header copied from paint_sgd_valid.gaf
!===============================
!Created on Mon Dec 7 11:33:04 2020.
!generated-by: PANTHER
!date-generated: 2020-12-07
!PANTHER version: v.15.0.
!GO version: 2020-11-17.
!
!===============================
!
!Documentation about this header can be found here: https://github.com/geneontology/go-site/blob/master/docs/gaf_validation.md
!
...
SGD S000004103 HOG1 GO:0003682 PMID:24508389 IDA F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0004707 PMID:10805732 IDA F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0005516 PMID:27421986 IPI UniProtKB:P06787 F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:10805732 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:12743037 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:23178807 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006972 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0007231 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0016241 PMID:16874103 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0033262 PMID:23178807 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0045944 PMID:12743037 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
...
```

# Downloading GO annotations

GAF format



```
!gaf-version: 2.1
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-09
!
!Header from source association file:
!==================================
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-08
!
!Header from sgd source association file:
!==================================
!Date: 20201207
!From: Saccharomyces Genome Database (SGD)
!URL: https://www.yeastgenome.org/
!Contact Email: sgd-helpdesk@lists.stanford.edu
!Funding: NHGRI at US NIH, grant number U41-HG001315
!
!==================================
!
!Header copied from paint_sgd_valid.gaf
!==================================
!Created on Mon Dec 7 11:33:04 2020.
!generated-by: PANTHER
!date-generated: 2020-12-07
!PANTHER version: v.15.0.
!GO version: 2020-11-17.
!
!==================================
!
!Documentation about this header can be found here: https://github.com/geneontology/go-site/blob/master/docs/gaf_validation.md
!
...
SGD S000004103 HOG1 GO:0003682 PMID:24508389 IDA F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0004707 PMID:10805732 IDA F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0005516 PMID:27421986 IPI UniProtKB:P06787 F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:10805732 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:12743037 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:23178807 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006972 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0007231 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0016241 PMID:16874103 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0033262 PMID:23178807 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0045944 PMID:12743037 IDA P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
...
```

**DB object symbol**

Physalia
Courses

# Downloading GO annotations

GAF format

# Downloading GO annotations

GAF format

```
!gaf-version: 2.1
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-09
!
!Header from source association file:
!=================================
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-08
!
!Header from sgd source association file:
!=================================
!Date: 20201207
!From: Saccharomyces Genome Database (SGD)
!URL: https://www.yeastgenome.org/
!Contact Email: sgd-helpdesk@lists.stanford.edu
!Funding: NHGRI at US NIH, grant number U41-HG001315
!
!=================================
!
!Header copied from paint_sgd_valid.gaf
!=================================
!Created on Mon Dec 7 11:33:04 2020.
!generated-by: PANTHER
!date-generated: 2020-12-07
!PANTHER version: v.15.0.
!GO version: 2020-11-17.
!
!=================================
!
!Documentation about this header can be found here: https://github.com/geneontology/go-site/blob/master/docs/gaf_validation.md
!
...
SGD  S000004103  HOG1  GO:0003682  PMID:24508389  IDA  F  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0004707  PMID:10805732  IDA  F  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0005516  PMID:27421986  IPI  UniProtKB:P06787  F  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0006468  PMID:10805732  IDA  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0006468  PMID:12743037  IDA  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0006468  PMID:23178807  IDA  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0006972  PMID:7681220  IMP  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0007231  PMID:7681220  IMP  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0016241  PMID:16874103  IMP  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0033262  PMID:23178807  IDA  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD  S000004103  HOG1  GO:0045944  PMID:12743037  IDA  P  Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
...
```

**Ref. for evidence**

# Downloading GO annotations

GAF format

```
!gaf-version: 2.1
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-09
!
!Header from source association file:
!=================================
!
!Generated by GO Central
!
!Date Generated by GOC: 2020-12-08
!
!Header from sgd source association file:
!=================================
!Date: 20201207
!From: Saccharomyces Genome Database (SGD)
!URL: https://www.yeastgenome.org/
!Contact Email: sgd-helpdesk@lists.stanford.edu
!Funding: NHGRI at US NIH, grant number U41-HG001315
!
!=================================
!
!Header copied from paint_sgd_valid.gaf
!=================================
!Created on Mon Dec 7 11:33:04 2020.
!generated-by: PANTHER
!date-generated: 2020-12-07
!PANTHER version: v.15.0.
!GO version: 2020-11-17.
!
!=================================
!
!Documentation about this header can be found here: https://github.com/geneontology/go-site/blob/master/docs/gaf_validation.md
!
...
SGD S000004103 HOG1 GO:0003682 PMID:2450838 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0004707 PMID:1080573 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0005516 PMID:2742198 IPI UniProtKB:P06787 F Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:1080573 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:1274303 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006468 PMID:2317880 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0006972 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0007231 PMID:7681220 IMP P Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0016241 PMID:1687410 IMP   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0033262 PMID:2317880 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
SGD S000004103 HOG1 GO:0045944 PMID:1274303 IDA   Mitogen-activated protein kinase involved in osmoregulation YLR113W|SSK3|mitogen-activated protein kinase HOG1
...
```

**Type of evidence**

When you have a **defined** set of <u>tens</u> to <u>hundreds</u> of genes

- o Genes significantly over-expressed (e.g. fold-change > 2 and adj. p-value < 0.01) in one condition vs a control

- o Genes whose promoter is bound by a specific combination of transcription factors

- o Thousands are probably too many genes…

# GO over-representation analyses: why

- Essentially, to know whether GO terms are over-enriched in a specific list of genes

- To get **an** idea of the functional/structural role of your set of genes

- To bring **a** piece of evidence that your treatment triggers some BP/MF/CC

- To know how much of the genes involved in a specific BP/MF/CC are present in your set of interest.

Finding over-represented GO terms in a given set of genes is one of the most common tasks in genomics.

Finding over-represented GO terms in a given set of genes is one of the most common tasks in genomics.

It usually relies on a straightforward Fisher test

Finding over-represented GO terms in a given set of genes is one of the most common tasks in genomics.

It usually relies on a straightforward Fisher test

| Guy who has drunk a ton of tea | Milk first | Tea added first |
|---|---|---|
| Cup is good | 23 | 52 |
| Cup is bad | 43 | 12 |



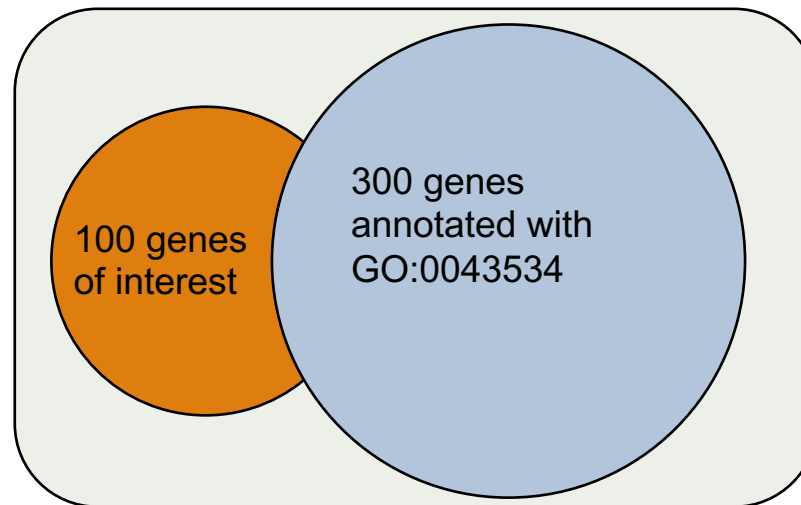Milk first / Tea first ?
Good / bad ?

Finding over-represented GO terms in a given set of genes is one of the most common tasks in genomics.

It usually relies on a straightforward Fisher test



100 genes of interest

300 genes annotated with GO:0043534

Universe: ~5,000 genes with at least 1 GO annotation

Finding over-represented GO terms in a given set of genes is one of the most common tasks in genomics.

It usually relies on a straightforward Fisher test

Think about it in terms of <u>contingency tables</u>

# GO over-representation analyses: how

| UNIVERSE = All Yeast genes annotated in the Biological Processes (5067 genes) | Genes over-expressed *in an assay* (152) | ~~Genes over-expressed *in an assay*~~ | |
|---|---|---|---|
| **Genes annotated** *in GO:0006836* *(243)* | 89 | 154 | → Sum = number of genes in GO:0006836 (243) |
| ~~**Genes annotated** *in GO:0006836*~~ | 63 | 5067 – 89 – 154 - 63 | |

↓
Sum = number of genes over-expressed in an assay (152)

Total sum = number of genes in BP (5067)

Physalia Courses

# GO over-representation analyses: how

| UNIVERSE = All Yeast genes annotated in the Biological Processes (5067 genes) | Genes over-expressed *in an assay* (152) | $\overline{\text{Genes over-expressed}}$ *in an assay* |
|---|---|---|
| Genes annotated *in GO:0006836* (243) | 89 | 154 |
| $\overline{\text{Genes annotated}}$ *in GO:0006836* | 63 | 5067 – 89 – 154 - 63 |

➔ Now repeat that for the 44,945 GO terms in the GO database……….

Physalia
Courses

# GO over-representation analyses: how

| UNIVERSE =<br>All Yeast genes annotated in the Biological Processes<br>(5067 genes) | Genes over-expressed<br>*in an assay*<br>*(152)* | ‾‾‾‾‾‾‾‾‾‾‾‾<br>Genes over-expressed<br>*in an assay* |
|---|---|---|
| Genes annotated<br>*in GO:0006836*<br>*(243)* | 89 | 154 |
| ‾‾‾‾‾‾‾‾‾‾‾‾<br>Genes annotated<br>*in GO:0006836* | 63 | 5067 – 89 – 154 - 63 |

➜ AND DON'T FORGET TO CORRECT FOR MULTIPLE TESTING (because testing 44,945 times is multiple testing…)

- Fortunately, there are many tools already out there to efficiently perform these calculations

- Some web-based, some with programmatic access

- They function with a range of "autonomy". Some need you to download the GO database, the GO annotations, or are doing all the work for you in the background

# Programs to run GO over-representation analyses: gProfiler

- Also available in R!

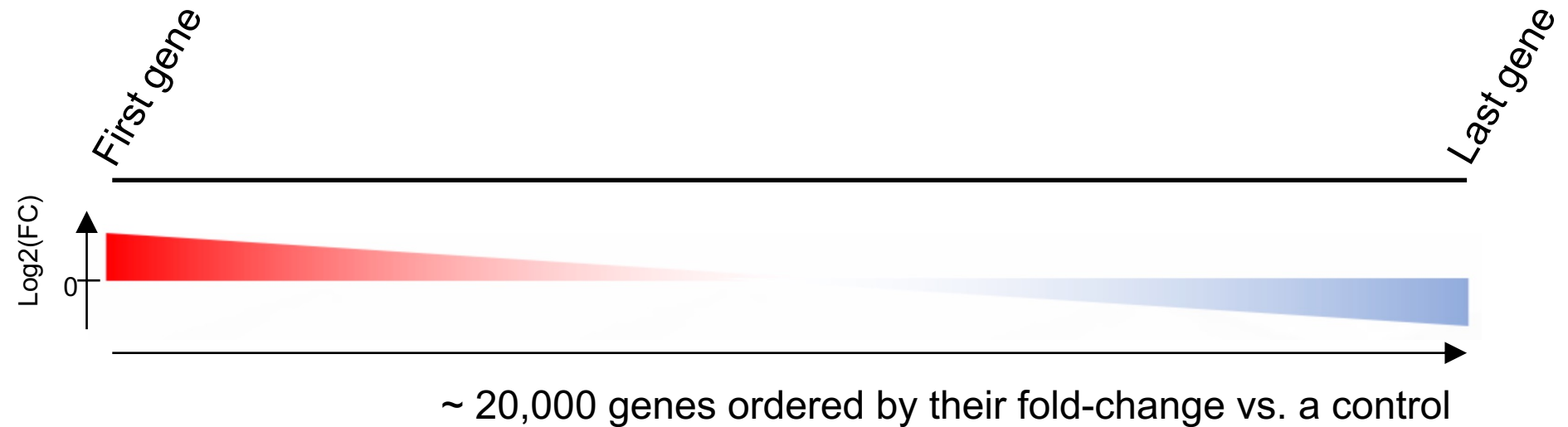- Simple, but many optional parameters to optimize your search

```
gprofiler2::gost(
    geneList,
    organism = 'scerevisiae’
)
```

# What if I don't have a gene __set__ of interest?

- Sometimes, you cannot really decide what is significant or not

- You don't like the idea of taking the top 100 genes differently expressed genes
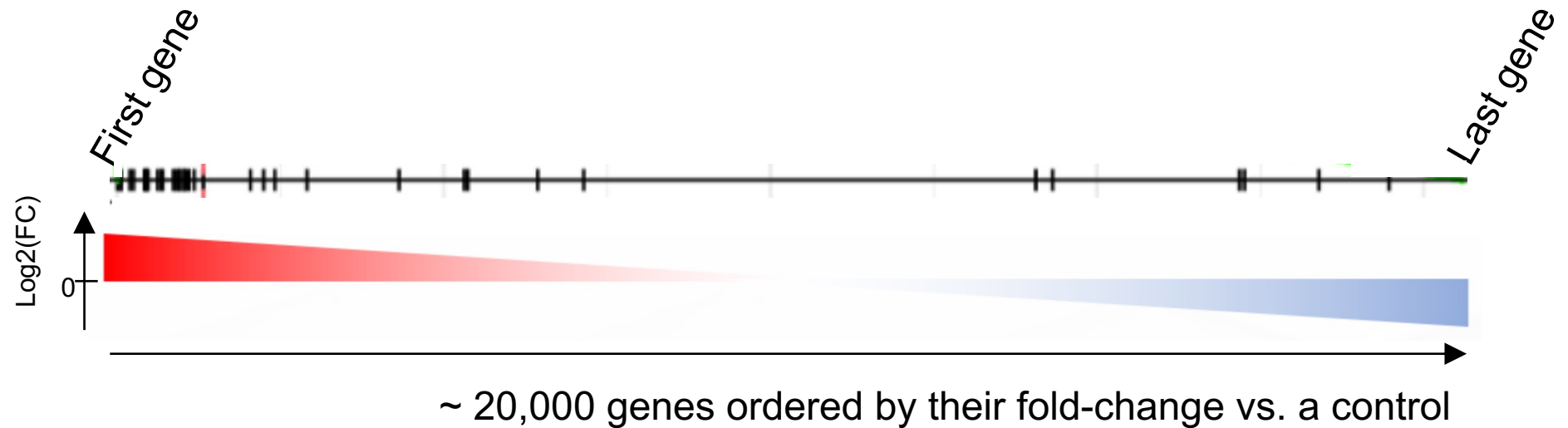
- How to set a threshold for your genes? FC>2? FC>5? ???

Physalia
Courses

GSEA (Gene Set Enrichment Analysis) uses a ranked list of genes as input



~ 20,000 genes ordered by their fold-change vs. a control

Within this list, it flags the genes belonging to a gene set (e.g. genes annotated in "centriole assembly" GO term)



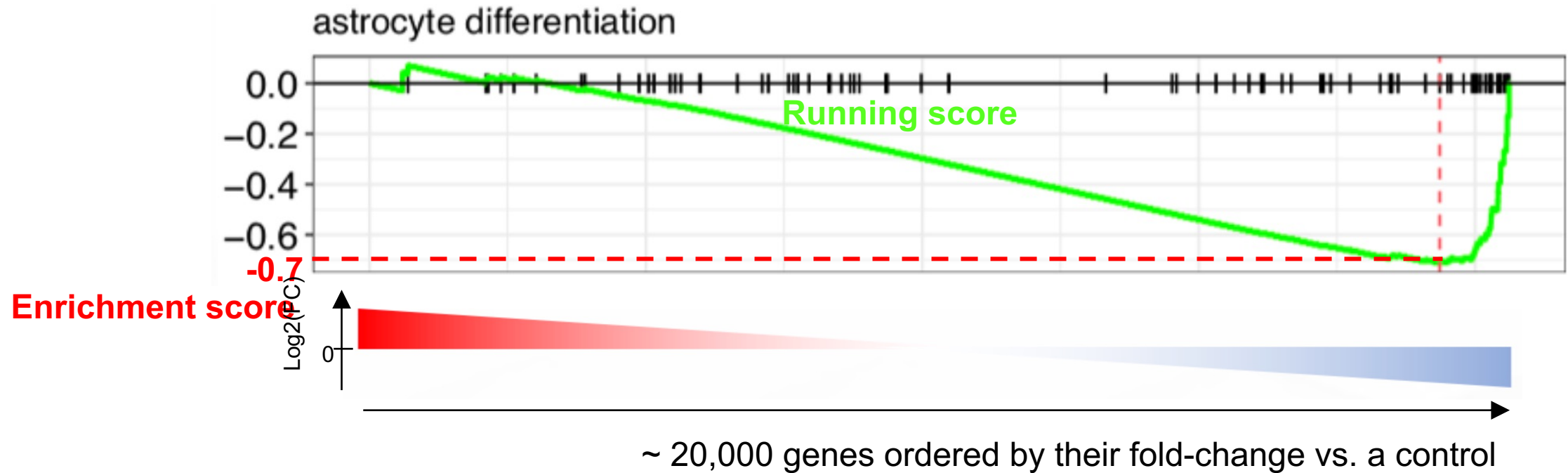~ 20,000 genes ordered by their fold-change vs. a control

# Gene Set Enrichment Analysis: cutoff-free functional analysis of ranked lists of genes

Based on the distribution of the flagged genes, it computes a "running score" and an "enrichment score"



~ 20,000 genes ordered by their fold-change vs. a control

# Gene Set Enrichment Analysis: cutoff-free functional analysis of ranked lists of genes

It can also find negative enrichment scores (indicated a depletion of genes of interest in the top of a ranked list)



~ 20,000 genes ordered by their fold-change vs. a control

# How to perform GSEA?

Original software: in JAVA

- o I never managed to use it…

# How to perform GSEA?
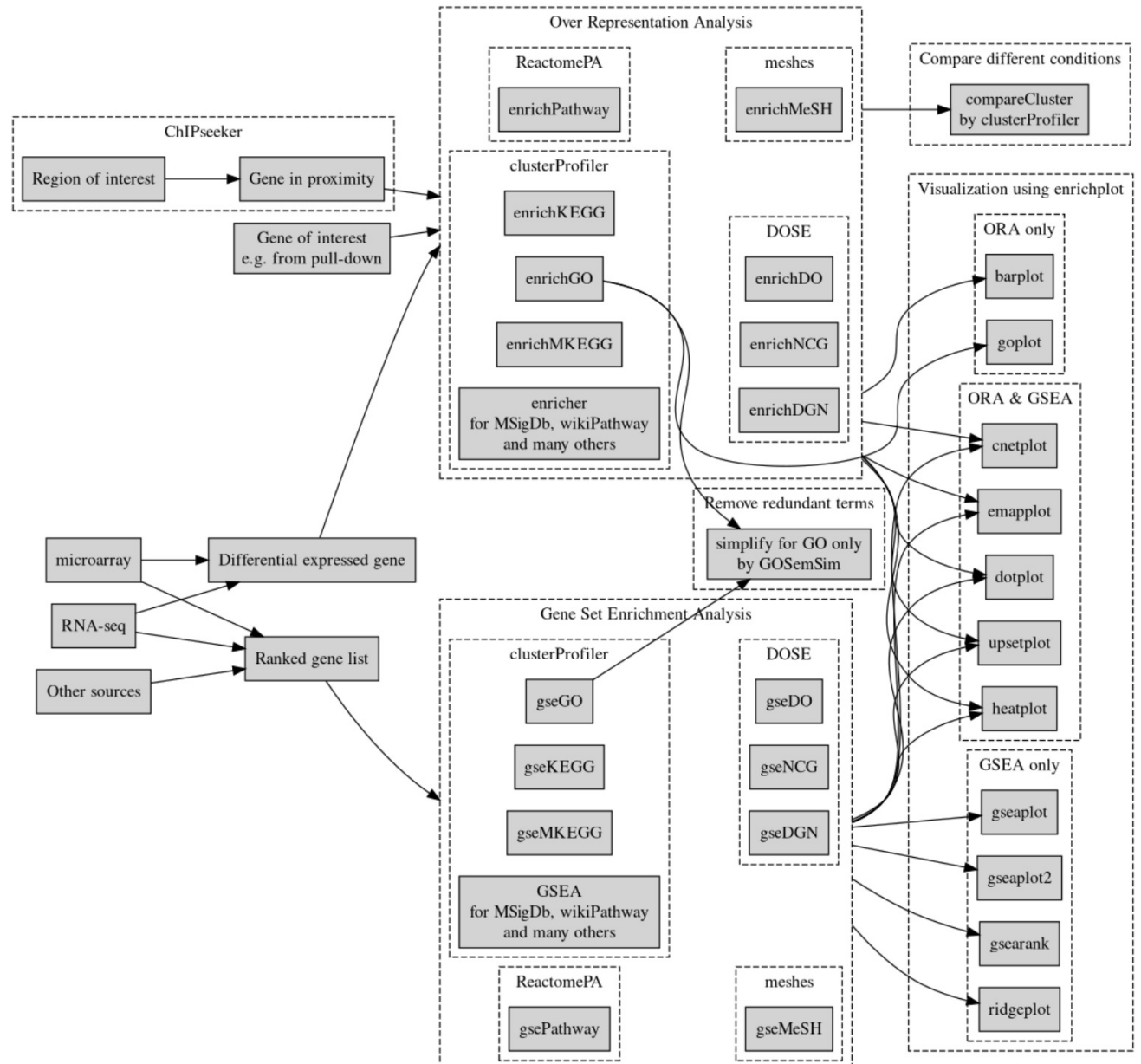
Original software: in JAVA

- o I never managed to use it…

Since then: many programmatic implementations notably in R

- o clusterProfiler is my personal favorite

- o Based on fgsea, the original GSEA implementation in R

- o Very complete and extensive doc

- o Nice visualization outputs

- o Well-integrated with GO ecosystem and other databases (disease ontology, Reactome, …)
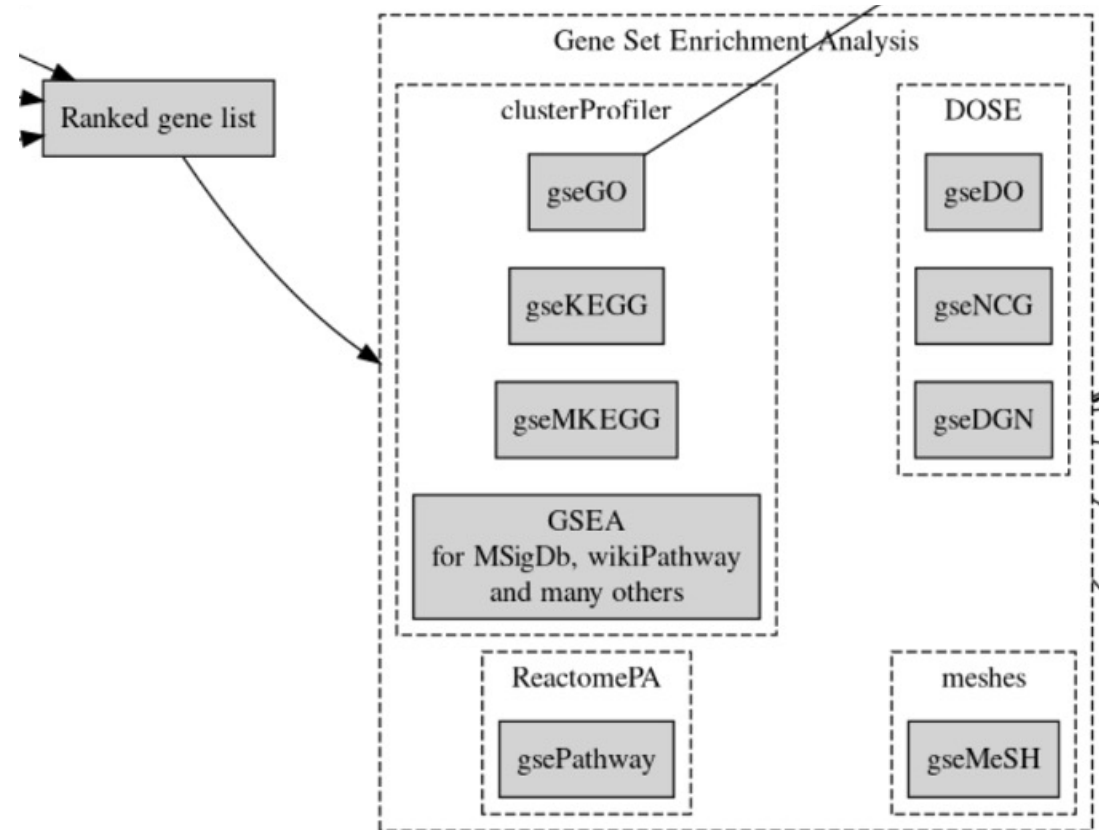
clusterProfiler is a rich set of tools to assess and visualize enrichment of a set of genes of interest compared to different databases

# clusterProfiler

clusterProfiler provides multiple gse*() functions, based on the type of gene sets you want to use
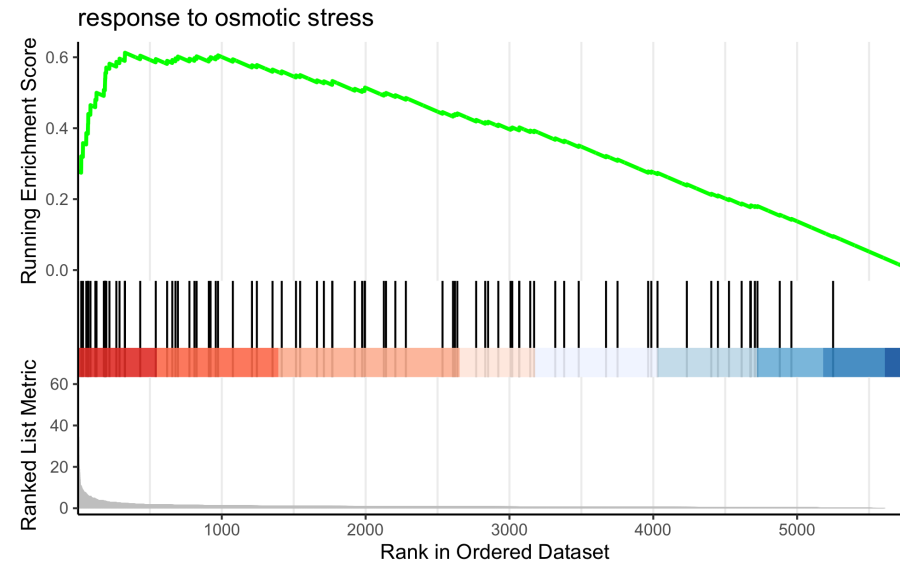
# clusterProfiler

- clusterProfiler provides multiple gse*() functions, based on the type of gene sets you want to use

- gseGO() compares your ranked list of genes to the reference up-to-date GO annotations

```
clusterProfiler::gseGO(
    rankedGeneLit,
    keyType = "ENSEMBL",
    OrgDb = org.Sc.sgd.db::org.Sc.sgd.db
)
```

# clusterProfiler

clusterProfiler provides plotting functions to check the results for a given GO term

```r
TERM <- "response to osmotic stress"
enrichplot::gseaplot2(
    gsea_results,
    title = TERM,
    geneSetID = which(gsea_results@result$Description == TERM)
)
```

*Ten Quick Tips for Using the Gene Ontology*, Blake PLoS Comp. Biol. 2013

# Tip 1a: Know the Source of the GO Annotations You Use
# Tip 1b: Know the Version of the GO Annotations You Use

*clusterProfiler: universal enrichment tool for functional and comparative study*, Guangchuang Yu
(http://yulab-smu.top/clusterProfiler-book/)