# Omics Playground Documentation

# GETTING STARTED

Omics Playground is a comprehensive self-service platform platform for visualization, analytics and exploration of Big Omics Data. It allows users to apply a multitude of state-of-the-art analysis tools to their own data to explore and discover underlying biology in a short time.

The platform offers a unique combination of features that distinguishes it from the other analytics platforms currently available. We believe that data preprocessing (primary analysis) and statistical testing (secondary analysis) are now well established, and the most challenging task is currently data interpretation (tertiary analysis) that often takes the longest time but where actual insights can be gained. Therefore Omics Playground focuses strongly on tertiary analysis while providing good support for secondary analysis.

# ONE

# OVERVIEW

Omics Playground is a user-friendly and interactive web-based platform for the analysis and visualization of transcriptomics and proteomics data. Currently the platform handles any kind of transcriptomics and proteomics data annotated at the gene/protein level, and supports two species, human and mouse. Omics Playground has been in particular devised to also support single cell RNA-seq data, as well as traditional gene expression experiments.

The overview of the platform is shown in the figure below. It consists of two main components. The first component addresses the data importing and preprocessing, which includes preparation of the input data, filtering, normalisation and precomputation of statistics for some analyses. The second part is composed of the online interface, which supports the real-time visualisation and interaction with users. The interface contains several optional settings in order to provide a customisable experience suited to each user's background.

The platform can be accessed through the following link: https://bigomics.ch/omics-playground/

## Data cleaning & preprocessing

| Quantification | Filtering | Normalisation | Offline computation | |
|---|---|---|---|---|
| Salmon, Kallisto, STAR or counts from any other software or database | Filtering based on the variance, abundance, or phenotype | Log2 transformation, quantile normalisation, batch effect correction | Precomputing the DE and Enrichment analysis to enable real-time visualisation | Offline |

## Omics-playground interface

### Home
Select and load the data

### Data table
Descriptive statistics

### Clustering
Heatmap    PCA/tSNE

### DE analysis
Volcano    Top DE genes

### Enrichment analysis
GSEA    GS activation

### Functional analysis
Pathway    Activation-map

### Intersection analysis
Enrichment & DE intersection    Cor. plots

### Signature analysis
Expression    Enrichment

### Cell profiling
Cell type prediction    Markers

Online, real-time visualisation

# INTRODUCTION

The platform requires the transcriptomics and proteomics data to be in a structured format as an input. The easiest way is to prepare two csv files: **counts.csv**, **samples.csv** and an optional **contrast.csv**. When these files are ready, users can upload them direcly in Omics Playground.

## 2.1 Input file requirements

1. *Counts file* Count or expression .csv file with gene on rows, samples as columns.

2. *Samples file* Samples .csv file with samples on rows, phenotypes as columns.

3. *Contrasts file* Contrast .csv file with samples on rows, conditions as columns. (optional)

**See also:**

If you have raw file formats, such as FASTQ files or LC-MS proteomics data (mzML, RAW, WIFF...), check our tutorials on how to prepare the counts matrix from these raw formats: *data preparation examples*.

# **SAMPLES FILE**

The samples file (*samples.csv*) contains the phenotypic information of each sample. The first column contains the sample name, which must be unique, and has to match the name given in the read counts file.

The samples file is a tabular (csv) file with the samples in the rows and the phenotypic data (metadata) in the columns. Note that the platform will not accept purely numerical values as phenotypes.

If we are analyzing a human study (it can be applied to any study) as seen in the `samples.csv` table below, the rows should be anonymized patients, identifyied uniquely by the first column (sample1, sample2. . . ), and the other columns would be sample metadata or phenotypes (hair color, country, weight, age, etc.).

|         | hair_color | country     | age   |
|---------|------------|-------------|-------|
| sample1 | blond      | Japan       | old   |
| sample2 | dark       | Switzerland | young |
| sample3 | blond      | USA         | young |
| sample4 | dark       | Switzerland | old   |
| sample5 | dark       | USA         | old   |

As mentioned above, the age was converted from numeric (12, 52, 87) to young and old, since the platform currently does not support continuous values.

---

**Note:** All phenotypes must contain at least one alphabet letter. This is done to avoid continuous values (as in the case of weight), since the platform expects discrete ranges. Having excessive numbers of phenotypic groups may also result in errors.

---

**See also:**

If you are familiar with R, you can think of the samples file as a data.frame object. We provide an example samples file that can be accessed by installing playbase `devtools::install_github("bigomics/playbase")` and running `playbase::SAMPLES`.

# COUNTS FILE

The file 'counts' contains the measurements (genes, proteins, etc..) for each sample listed in the samples file. Just like the samples, the `counts.csv` file is tabular (.csv), where each row describes the features (genes, proteins, etc..) and each column describes the samples.

The rows contains gene IDs, which can be in most common formats (such as HGCN or Ensembl), but not in the Entrez number format. If you are using Entrez numbers, please convert them to Ensembl IDs using tools such as Syngo.

The values should always be numerical, with the exception of "NA" in case of a lack of data. Failure to do so will result in an error.

Below is a simple example of how a `counts.csv` file should look like.

|  | sample1 | sample2 | sample3 | sample4 | sample5 |
|---|---|---|---|---|---|
| gene1 | 543.6 | 1556.1 | 413.0 | 887.9 | 123.4 |
| gene2 | 6.5 | 14.7 | 2.3 | 42.4 | 56.7 |
| gene3 | 10.4 | 763.5 | NA | 0 | 89.0 |
| gene4 | 3217.4 | 0 | 4983.2 | 7493.8 | 210.2 |
| gene5 | 98770.5 | 113498.0 | 498351.6 | 88134.1 | 345.6 |
| gene6 | 0 | NA | 14.9 | 0 | 789.0 |
| gene7 | 47648.8 | 0 | 32682.0 | 93873.2 | 123.4 |

**Note:** The formats accepted as features (genes, proteins are ENSEMBL, ENSEMBLTRAN, UNIGENE, REFSEQ, ACCNUM and UNIPROT and gene SYMBOL). Also note that the platform will not accept transcript IDs. You will need to convert them to Gene IDs. This will result in multiple gene entries that the platform will merge.

**See also:**

If you are familiar with R, you can think of the counts file as a data.frame object. We provide an example samples file that can be accessed by installing playbase `devtools::install_github("bigomics/playbase")` and running `playbase::COUNTS`.

# CONTRASTS FILE (OPTIONAL)

The contrasts are an optional input in Omics Playground. It contains the groups (treatment versus controls, mutant versus wild-type, etc..) that will be tested against each other. The contrasts are a list of comparisons between groups.

There are two ways to define the contrasts: the long contrast form and the short contrast form. The contrast name rows need to be unique, both in short and long forms.

## 5.1 Sample-wise contrast (long form)

The long contrast form has the same number of rows of the samples file. That means the first column is the sample names (matching samples file), and the second column (onwards) contains the phenotypes.

The first row contains the name of the pairwise comparisons. All pairwise contrasts names must follow the format shown below with the groups joined together by "_vs_" (e.g. piperaquine_vs_control).

**..note::**

1. All comparison names in the header **must** contain the separator "_vs_"

2. Please use the exact condition names in the title as well as in their columns

For example, if you have 3 groups (A, B, C), you can test the following contrasts: A_vs_B, A_vs_C, B_vs_C.

Here is a minimal example of how the `contrasts.csv` should look like. In this case, the groups will be hair color (blond vs. dark) and country (Japan vs. Switzerland), as derived from the samples file. Use NA (or empty field) to skip samples that do not belong to a contrast.

|  | dark_vs_blond | japan_vs_switzerland |
| --- | --- | --- |
| sample1 | blond | japan |
| sample2 | dark | switzerland |
| sample3 | blond | japan |
| sample4 | dark | switzerland |
| sample5 | dark | NA |

The column names (dark_vs_blond) will be used to name the comparisons. The first name (dark) will be used as the numerator of the contrast, and the second name (blond) will be used as the denominator of the contrast. Alternatively, you can code the groups with -1 (reference group), 1 (main group) and 0 (not in comparison). This is shorter than writing the entire word but you have to be careful which group to assign +1 and -1 to. The -1 refers to the 'control' group that is the denominator or 'on-the-right' of the '_vs_' in the comparison name.

| | dark_vs_blond | japan_vs_switzerland |
|---|---|---|
| sample1 | -1 | 1 |
| sample2 | 1 | -1 |
| sample3 | -1 | 1 |
| sample4 | 1 | -1 |
| sample5 | 1 | 0 |

Any of the tables above can be provided to the platform. The zero means the sample5 will not be include in the comparison of japan_vs_switzerland.

## 5.2 Group-wise contrasts (short form)

Sample-wise contrasts, as introduced in the previous paragraph, become tedious in datasets with large number of samples, as assigning the contrasts to each sample can be cumbersome. Instead, you can specify group-wise contrasts (or short form). However, this approach only works best if users are focusing on a single phenotype in their dataset.

Following the example above, if we would like to create a contrast between the two countries Japan and Switzerland, we will need to create a column called **group** in the sample file, which will contain the phenotypes (in this case country). To simplify, we could simply change the name of the column **country** to **group**.

Once we have the group column in the sample file, we can assign the contrasts as in the table below.

| | japan_vs_switzerland | USA_vs_switzerland | all_vs_switzerland |
|---|---|---|---|
| japan | 1 | 0 | 1 |
| USA | 0 | 1 | 1 |
| switzerland | -1 | -1 | -1 |

We will search for the group column in the samples file, and we will create the contrasts based on the groups. For the first comparison, samples with 'japan' in the group column will be the numerator, and samples with 'switzerland' in the group column will be the denominator. For the second comparison sample with 'USA' (as main group) will be compared to the samples with 'switzerland' (as reference group). The third comparison specifies both Japan and USA (together as one group) to be compared to Switzerland.

**..note::**

1. Group-wise contrasts require a **group** column in the samples file.

2. The group column name must contain the word *group* (e.g. group, groups, group_name, etc..).

3. If multiple group columns are found, only the first one will be matched.

**..seealso::**

If you are familiar with R, you can think of the contrasts file as a data.frame object. We provide an example samples file that can be accessed by installing playbase `devtools::install_github("bigomics/playbase")` and running `playbase::CONTRASTS`.

# CHECKS BEFORE UPLOADING

**Check #1**. Never use the special characters in your samples and contrasts, as certain characters are reserved for programming purposes. If you need to connect multiple elements, use underscore, "_" instead.

We test regularly which characters are allowed or not, we will update this list accordingly. The characters below are strictly forbidden.

| Input file | Part | Characters not allowed |
|---|---|---|
| sample.csv | column names | @ , : empty spaces + |
| sample.csv | row names | @ , : empty spaces |
| sample.csv | phenotype names | @ , : |
| contrasts.csv | column names | @ , : empty spaces + |
| contrasts.csv | contrast names | @ , : |

**Note:** We cannot guarantee that special characters (+, -, *, /, %, etc..) and empty space will work in all Omics Playground modules, so we recommend substituting them with underscore '_'.

**Check #2** Avoid starting your sample, phenotype or contrast names with special characters like (+, -, *, /, %, etc..). While the platform will accept these characters, they may be converted into a standard symbol like X. For example, if we label the sample.csv country column as *%country*, we will see X.country in some analysis in Omics Playground.

**Check #3** Make sure you have the words 'samples', 'counts' and 'contrasts' in the corresponding filenames. For example, we accept experiment23_samples.csv, samples_experiment23.csv, but not experiment23.csv as sample input.

**Check #4**. Define intervals instead of using numeric phenotypes. The platform does not cope with continuous numeric variables for phenotypes yet. To avoid that, our coders added a filter that flags phenotypes names "Time" or "Age" as unacceptable. The same applies for other continuous variables, such as height, weight, length, etc.

Instead, you should cluster the various numeric values into definite intervals and then name them accordingly (e.g. "Age_groups", "Time_intervals", etc...)

|          | Age |
|----------|-----|
| Patient1 | 28  |
| Patient2 | 35  |
| Patient3 | 37  |
| Patient4 | 72  |
| Patient5 | 45  |

|          | Age_group |
|----------|-----------|
| Patient1 | 26_35     |
| Patient2 | 25_35     |
| Patient3 | 36_45     |
| Patient4 | over70    |
| Patient5 | 36_45     |

# UPLOADING YOUR DATA IN OMICS PLAYGROUND

Users can import their data from the **Upload data** panel located under the *Load Panel* module. The platform requires a file with the read counts and one with a description of the samples at the minimum. An optional file with the desired contrasts can also be provided. The format of files must be comma-separated-values (CSV) text. It is important to name and format the files as explained in the previous sections.

# WHAT IF I HAVE MANY DATASETS?

The first step in any OmicsPlayground analysis is to upload data which can be used to create a pgx object. The pgx object is basically the core data structure in the OmicsPlayground upon which most analysis and plotting functions operate.

If you have many datasets (tens to hundreds), it might make sense to create the pgx object with an R script.

We made it very easy to create a pgx object from your own data. With a few R functions, you can convert several datasets into pgx objects that can be uploaded in Omics Playground.

Here we check that your input files do not have problems

```r
# install necessary packages

install.packages("devtools")
devtools::install_github("bigomics/playbase")

library(playbase)

# These are the possible errors you can encounter

playbase::PGX_CHECKS

# individual file checks

SAMPLES = playbase::pgx.checkINPUT(playbase::SAMPLES, type = "SAMPLES")
COUNTS = playbase::pgx.checkINPUT(playbase::COUNTS, type = "COUNTS")
CONTRASTS = playbase::pgx.checkINPUT(playbase::SAMPLES, type = "CONTRASTS")

# Checks across input files

INPUTS_CHECKED <- pgx.crosscheckINPUT(SAMPLES, COUNTS, CONTRASTS)

SAMPLES = INPUTS_CHECKED$SAMPLES
COUNTS = INPUTS_CHECKED$COUNTS
CONTRASTS = INPUTS_CHECKED$CONTRASTS
```

If no errors are reported (and *PASS* is *TRUE*), these new checked files SAMPLES, COUNTS and CONTRASTS can be used safely in the next step.

Here we create a pgx object that can be used in Omics Playground:

```r
# step 1: create a pgx object with your samples, counts and contrasts
pgx <- playbase::pgx.createPGX(
```

```
  counts = playbase::COUNTS,
  samples = playbase::SAMPLES,
  contrasts = playbase::CONTRASTS
)

# step 2: compute the pgx object

pgx <- playbase::pgx.computePGX(
  pgx = pgx
)

# save the pgx object to your computer
save(pgx, file = "choose_a_name.pgx.")
```

All you have to do is substitute `playbase::COUNTS`, `playbase::SAMPLES`, and `playbase::CONTRASTS` with your own data. You can then import the pgx object into Omics Playground.

**See also:**

In reality, there is a lot more happening behing the *pgx.createPGX* and *pgx.computePGX*. If you are interested in learning more, please see our Github Wiki with more details on the statistics, normalization and filtering steps that are performed. Wiki. You can find it here.

# NINE

# PREPARING INPUT FILES FROM RAW DATA FORMATS

If you have a raw data format (FASTQ, mzML, mzML, RAW, WIFF, etc..) or want to perform additional pre-processing steps (normalization or filtering), in this section we show tutorials on how to perform such steps. We also show how to prepare data files from public databases such as GEO.

We provide four types of example cases to guide users for preparing their input objects and injecting it into the platform. Basically, the example cases illustrate how to prepare an input data:

1. *from FASTQ files*,

2. *from gene counts table or from the GEO repository*,

3. *from single-cell data*,

4. *from LC-MS/MS proteomics data*.

All the necessary scripts for data cleaning and preprocessing examples can be found under the `scripts/` folder.

## 9.1 From FASTQ files

Starting from FASTQ files, we recommend using the GREP2 package to obtain gene counts through quality control, trimming, quantification of gene abundance, and so on. Afterwards, the user can refer to the examples in the next section for preparing an input data from the gene counts.

## 9.2 From gene counts table or GEO repository

Users can prepare an input data from their own gene counts or download a relevant dataset from repositories such as GEO. Some examples are provided in the following scripts:

- TCGA-BRCA: `pgx-tcga-brca.R`

- TCGA-PRAD: `pgx-tcga-prad.R`

- GSE10846: `pgx-GSE10846-dlbcl.R`

- GSE114716: `pgx-GSE114716-ipilimumab.R`

- GSE22886: `pgx-GSE22886-immune.R`

- GSE28492: `pgx-GSE28492-roche.R`

- GSE32591: `pgx-GSE32591-lupusnephritis.R`

- GSE53784: `pgx-GSE53784-wnvjev.R`

- GSE88808: `pgx-GSE88808-prostate.R`

## 9.3 From single-cell data

Single-cell RNA sequencing experiments have been valuable to provide insights into complex biological systems, reveal complex and rare cell populations, uncover relationships between genes, and track the trajectories of cell lineages. Below we provide some data preparation examples from single-cell experiments:

- GSE72056: `pgx-GSE72056-scmelanoma.R`

- GSE92332: `pgx-GSE92332-scintestine.R`

- GSE98638: `pgx-GSE98638-scliver.R`

## 9.4 From LC-MS/MS proteomics data

Two examples are provided below for LC-MS/MS proteomics data preprocessing:

- Geiger et al. 2016: `pgx-geiger2016-arginine.R`

- Rieckmann et al. 2017: `pgx-rieckmann2017-immprot.R`

# LOADING CUSTOM GENESET FILES

GMT (Gene Matrix Transposed) files are a format commonly used in bioinformatics and gene set enrichment analysis (GSEA).

Gene sets are collections of genes that share certain biological characteristics, such as being involved in the same biological pathway or having similar functions. GSEA is a computational method used to determine whether a particular gene set is significantly enriched in a given dataset of gene expression.

In a GMT file, the first column is the geneset name, followed by a description on the second column. Each subsequent column line represents a gene set and consists of a gene set name, and the list of genes belonging to that set. The gene names are typically represented using gene symbols or other identifiers. Each line in a GMT file will correspont to a different geneset.

Download here an example GMT file from genes activated by the transcription factor EGFR: `EGFR_TARGET_GENES.v2023.1.Hs.gmt`

---

**Tip:** If you want to create your own genesets, the easiest way is to download our example GMT file and open it in excel as tab-separated values. The first column is the geneset name (EGFR_TARGET_GENES), followed by any description (can be left empty) or the geneset source website in the second column. Then finally the genes should be listed on the third and following columns (ABCA7, ACTB, etc.). After editing, save the file as a tab-separated values file and change the extension to .gmt.

---

In Omics playground, GMT files are the backbone of several modules, where we aim to identify enriched gene sets associated with specific biological processes, diseases, or experimental conditions.

You can find the geneset availables in omicsplayground with the R package playbase. First, install playbase with `devtools::install_github("omicsplayground/playdata")` and then run the command `playdata::GSET_INFO`.

If you have your own geneset, or you want to include a GMT file from a database, you can upload directly in the Upload module under Options.

**See also:**

If you are working in R, you can add your gmt file to the `playbase::pgx.computePGX` by providing the `custom.geneset` argument. custom.geneset argument. That should be a list where the first argument `gmt` is the gmt file (see `playbase::EXAMPLE_GMT`), and the second argument is the name of the gmt file.

Building on what we learned in the section *Computing PGX* , we can add our own geneset to the analysis with the followign call:

```
filePath <- getwd() # path to gmt file

custom.geneset <- list()
```

```r
custom.geneset$gmt <- playbase::read.gmt(filePath)

# perform some basic checks
gmt.length <- length(custom.geneset$gmt)
gmt.is.list <- is.list(custom.geneset$gmt)

# clean genesets

custom.geneset$gmt <- list(CUSTOM = custom.geneset$gmt)

# convert gmt to OPG standard
custom.geneset$gmt <- playbase::clean_gmt(custom.geneset$gmt,"CUSTOM")

# compute custom geneset stats
custom.geneset$gmt <- custom.geneset$gmt[!duplicated(names(custom.geneset$gmt))]
custom.geneset$info$GSET_SIZE <- sapply(custom.geneset$gmt,length)

# pass the custom.geneset as argument to pgx.computePGX
pgx <- playbase::pgx.computePGX(
    pgx = pgx,
    custom.geneset = custom.geneset
)
```

# ELEVEN

## DATA PROCESSING INSIDE OMICS PLAYGROUND

## 11.1 Filtering of features (genes) and samples

The data preprocessing includes some filtering criteria, such as filtering of genes based on variance, the expression across the samples, and the number of missing values. Similarly, samples can also be filtered based on the read quality, total abundance, unrelated phenotype, or an outlier criterion.

## 11.2 Normalisation

The raw counts are converted into counts per million (CPM) and log2. Depending on the data set, a quantile normalization can be applied. Known batches in the data can be corrected with limma or ComBat. Other unknown batch effects and unwanted variation can be further removed using surrogate variable analysis in the sva package.

## 11.3 Offline computation

Statistics for the differentially expressed genes analysis and gene set enrichment analysis are precomputed to accelerate the visualisation on the interface.

OUTLINE

## 12.1 Main menu

Using the main menu on top, you can navigate through the different analysis modules. Generally you want to start from top to bottom, from specific gene-wise to the more higher-level functional analysis modules. Some users prefer the other way around.

By default, the following menu items are present, namely Load, *Data View*, *Clustering*, *Expression*, *GeneSets*, *Compare* and *SystemsBio*.

## 12.2 Help menu

The **Help** menu is located on the top right hand corner. The menu links to the online platform documentation (the page you are currently viewing), a collection of video tutorials on youtube, the community forum of the platform, the Github webpage, where busg can be reported, and a collection of case studies based on Omics Playground.



## 12.3 User menu

The **User** menu is located on the top right hand corner, right next to the **Help** menu. It links to the user profile, where information on the usage of the platform can be found, the app settings, the "About" pop-up that shows information about the version of the platform currently deployed and finally the "Logout" option to disconnect from the platform.



Clicking "App settings" will take the user to a new page with two tabs: **App settings & News** and **Resource info**.

There are two panels under the **App settings & News** tab: **Application options** and **New features**.

**Application options**

Through this panel, users can enable beta features, disable alerts and enable captions for the plots.

**New features**

This panel provides a list of new features implemented with each new update of the platform.

Three panels are found under **Resource info**: **Timings**, **PGX slot sizes** and **R object sizes**.

**Timings**

The timings table reports more detailed information about the object dimensions, object sizes and execution times of the methods.

**PGX slot sizes**

This table provides details about the pgx object.

**R object sizes**

This table provides size details about R objects.

## 12.4 Figure & table tags

Each figure or table on the platform is assigned the following interactive buttons, where:

- Info: provides detailed information about the figure or table.

- Settings: users can specify additional settings if applicable.

- Download: downloads a figure as a PNG or PDF file or a table in CSV format.

- Maximize: shows a larger version of a figure in a separate window.



## 12.5 Glossary

- Signature: a list of selected genes (e.g. by significance or fold change),

- Condition: a specific phenotype group (e.g. tumor or control),

- Contrast: a comparison between two conditions (e.g. tumor vs control),

- Profile: a vector of fold changes corresponding to a certain comparison,

- Hierarchical clustering: a method that groups similar samples into groups,

- Q value: an FDR-adjusted p value,

- Biomarker: a biological feature (gene, mutation or gene set) that characterises a specific physiological or pathological process.

# LOAD PANEL

## 13.1 Selecting a dataset

The **Load** panel panel shows the available datasets within the platform. Each dataset contains a brief description as well as the total number of samples, genes, corresponding phenotypes and the collection date.

Users can select a dataset in the table and then click on the "Load dataset" button at the bottom left to load a dataset into the platform. On the right hand side a signature t-SNE plot is displayed . Each dot corresponds to a specific comparison. Signatures/datasets that are clustered closer together, are more similar.



A pop-up menu is displayed when clicking on the three dots next toa dataset name. From this menu, users can download either a pgx object, which contains all the calculations and plots, or the raw input data that was used to generate the dataset. The "Reanalyse" option is used to perform a new analysis based on the same input files as the selected dataset. Users can also share a dataset with another user ( via the "Share with user" option) or make a dataset public (via the "Share public" option). Finally, "Delete dataset" is used to remove a dataset permanently from the platform.

N.B. Some of these features are (such as the delete option) are not available on trial accounts.

## 13.2 Load (Upload new)

Under the **Upload new** option of the **Load** panel users can upload their transcriptomics and proteomics data to the platform. The platform requires three data files as listed below: a data file containing counts/expression (counts.csv), a sample information file (samples.csv) and an optional file specifying the statistical comparisons as contrasts (contrasts.csv). The file format must be comma-separated-values (CSV) text. Be sure the gene names match for all files. On the left side of the panel, users need to provide a unique name and brief description for the dataset while uploading.

**counts.csv**
Count/expression file with gene on rows, samples as columns.

**samples.csv**
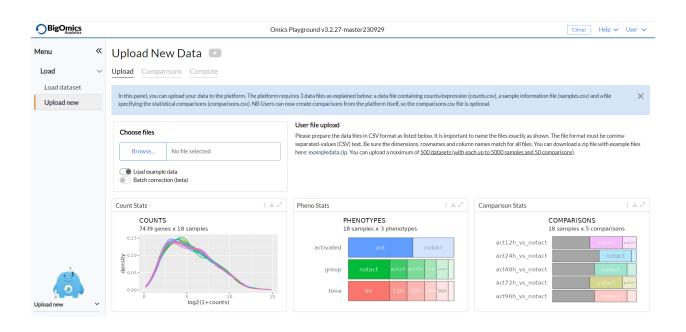Samples file with samples on rows, phenotypes as columns.

**contrasts.csv**
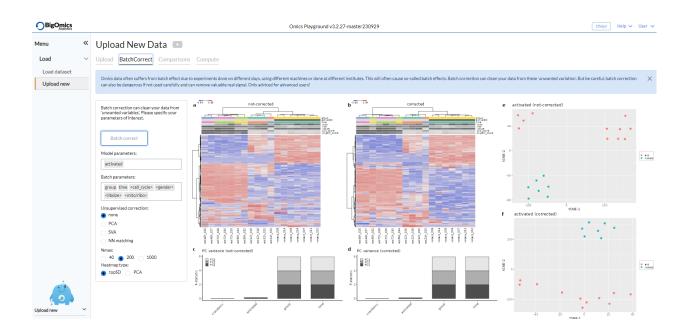Contrast file with conditions on rows, contrasts as columns.

Once uploaded, the platform generates three separate plots showing some stastics and the structure of your input files.

The page also contains a *batch correction* option that allows users to perform batch correction on the data. Selecting it will open a new panel from which users can select the level of batch correction (low, medium or strong). Under *Advanced* users can fine-tune the process. This feature is still in beta and only recommended to users familiar with the batch correction precedure and parameters. More information can be found under *Methods*.
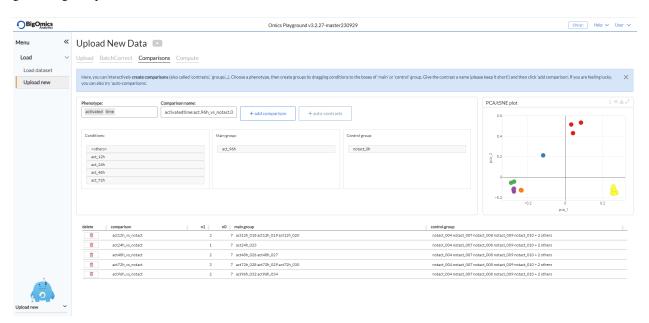
## 13.3 Upload new (Comparisons)

Users can eschew adding a contrast file and instead generate contrasts via the **Comparisons** subpanel of the **Upload new** panel. Available phenotypes will appear in the "Phenotypes" box as a scrolldown menu. Selecting a phenotype will show the available conditions in the "Conditions" box. Users can then drag individual groups in the "Main group" or "Control Group" boxes. It is also possible to select individual samples from the "Conditions" box and phenotypes can be furthermore combined in the "Phenotypes" box. The platform automatically generates a contrast name in the "Comparison name" box that users can manually edit. Clicking the "add comparison" button will add the selected contrast in the contrast table below. Users can also task the platform to generate comparisons automatically by selecting a phenotype form the "Strata" box and clicking on "add auto-contrasts". However, we recommend caution when using

this function as it is error-prone. The panel also produces a PCA/tSNE plot that users can consult as a guide for generating comparisons.



## 13.4 Upload new (Compute)

The **Compute** subpanel is where users can start the computation of their data. They need to provide a unique name for the dataset, indicate the type of data and provide a short description. Beginners can then click on the "Compute!" button and let the platform run the calculations. Advanced users can click on "Advanced" to access several customisation options. In particular they can de-select feature filtering options, select different gene tests and enrichment methods combinations, include or exclude analysis types and select developer-specific options. A new advanced feature with V3 of the platform is the ability to upload a custom GMT files with a dataset.

# FOURTEEN

# DATA VIEW

The **DataView** module provides information and visualisations of the dataset to quickly lookup a gene, check the counts, or view the data tables. It has five panels, which are briefly explained below, followed by more detailed information for each panel.

The **Gene Overview** panel displays figures related to the expression level of the selected gene, correlation, and average expression ranking within the dataset. In the **Sample QC** panel, the total number of counts (abundance) per sample and their distribution among the samples are displayed. This is most useful to check the technical quality of the dataset, such as total read counts or abundance of ribosomal genes. In **Counts table** panel, the exact expression values across the samples can be looked up, where genes are ordered by the correlation with respect to the first gene. Gene-wise average expression of a phenotype sample grouping is also presented in this table. In the **Sample information** panel, more complete information about samples can be found. Finally, the **Contrasts** panel, shows information about the phenotype comparisons.

## 14.1 Settings panel

Users can find the settings panel on the right hand side. The panel contains the main settings for the analysis. The analysis can be started by selecting a gene of interest from the `Gene` settings. Users can filter and select samples in the `Filter samples` settings, or collapse the samples by predetermined groups in the `Group by` settings. Under *Options*, it is possible to visualize the information on a raw count level or logarithmic expression level (logCPM).

## 14.2 Gene Overview

The **Gene Overview** panel displays figures related to the expression level of the selected gene, correlation to other genes, and average expression ranking within the dataset. To find out more information from the literature, hyperlinks are provide to connect the selected gene to OMIM, KEGG, and GO databases. It also shows the correlation of the gene in other datasets such as ImmProt and HPA, and plots the cumulative correlation. Furthermore, it displays the tissue expression for a selected gene using the genotype-tissue expression GTEx dataset. For each chart of the panel (left to right, top to bottom), a detailed explanation is provided below.

**Gene info**

To find out more information from the literature, hyperlinks are provide to connect the selected gene to public databases, including OMIM, KEGG, and GO.

**Gene Expression**

Expression barplot of grouped samples (or cells) for the selected gene. Samples (or cells) in the barplot can be ungrouped by setting the `Group by` under the main *Options*.

**Average Rank**

Ranking of the average expression of the selected gene.

**t-SNE clustering**
> T-SNE clustering of samples (or cells) colored by an expression of the gene selected in the `Search gene` dropdown menu. The red color represents an over-expression of the selected gene across samples (or cells).

**Top correlated genes**
> Barplot of the top positively and negatively correlated genes with the selected gene. Absolute expression levels of genes are colored in the barplot, where the low and high expressions range between the light blue and dark blue colors, respectively.

**Tissue expression (GTEX)**
> Tissue expression for the selected gene in the tissue expression GTEx database. Colors corresponds to "tissue clusters" as computed by unsupervised clustering.

## 14.3 Sample QC

In the **Sample QC** panel, the total number of counts (abundance) per sample and their distribution among the samples are displayed. For each sample, users can also see the percentage of counts in terms of major gene types such as ribosomal genes, heatshock proteins, or kinases. Abnormal abundance of certain genes may indicate technical problems. A detailed explanation is provided below for every chart of the panel (left to right, top to bottom).

**Total counts**
> A barplot of the total number of counts (abundance) for each group. The samples (or cells) can be grouped/ungrouped in the `Group by` setting uder the main *Options*.

**Median counts distribution**
> A boxplot of the total number of counts (abundance) for each group.

**Density distribution of counts**
> A plot showing the density distribution of counts for each group.

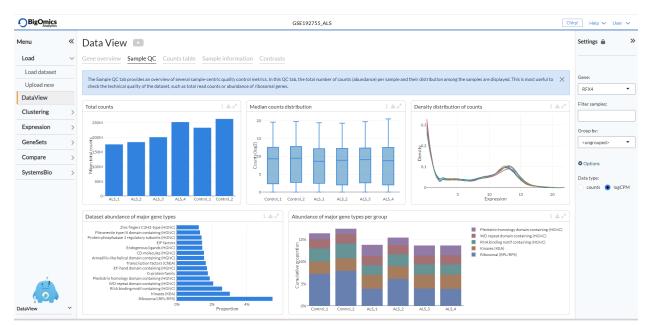**Dataset abundance of major gene types**
> A barplot showing the percentage of counts in terms of major gene types such as CD molecules,

kinanses or RNA binding motifs for each group.

**Abundance of major gene types per group**

A barplot showing the average count levels of major gene types such as CD molecules, kinanses or RNA binding motifs for each group.

## 14.4 Counts Table

Under the **Counts table** panel, the exact expression values across the samples can be read, where genes are ordered by the correlation with respect to the first gene. Gene-wise average expression of a phenotype sample grouping is also presented in this table.

The samples (or cells) in the table can be ungrouped by setting the `Group by` under the main *Options* to see the exact expression values per sample (or cell). The genes in the table are ordered by the correlation (**rho** column) with respect to the selected gene. **SD** column reports the standard deviation of expression across *all* samples (or cells).



## 14.5 Sample Information

In the **Sample information** panel, users can check information about samples and their phenotype grouping through three outputs (left to right, top to bottom):

**Phenotype clustering**
A plot showing phenotype clustering. Phenotypes can be unclustered via the *Settings* icon

**Phenotype association**
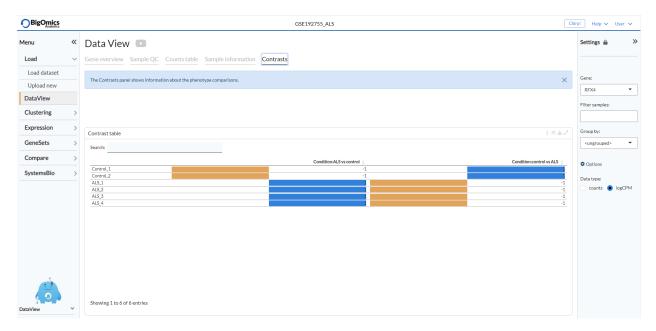A phenotype association matrix showing correlation between phenotypes.

**Sample information**
A table with sample information.

## 14.6 Contrasts

In the **Contrasts** panel, more complete information about contrasts can be found. It summarizes the contrasts of all comparisons. Here, users can check which samples belong to which groups for the different comparisons. Non-zero entries '+1' and '-1' correspond to the group of interest and control group, respectively. Zero or empty entries denote samples not use for that comparison.



Through the settings icon on top of the plot, users can display samples individually or in groups.

Show by:

○ sample

● group

# **CLUSTERING**

The **Clustering** module performs unsupervised clustering analysis of the data. After having done the QC, it is probably the first way to explore your data. The main purpose is to discover patterns and subgroups in the data, show correlation with known phenotypes, detect outliers, or investigate batch effects.

The module is divided into two submodules: **Samples** and **Features**.

Under the **Samples** submodule you can find classical clustering functions. In the **Heatmap** panel hierarchical clustering can be performed on gene level or gene set level. The **PCA/tSNE** panel shows unsupervised clustering of the samples in 2D/3D as obtained by PCA or tSNE algorithms. The **Parallel** panel displays the expression levels of selected genes across all conditions. On the right, the **Annotate cluster** panel provides a functional annotation for each feature cluster in the heatmap. Users can select from a variety of annotation databases from the literature. The **Phenotypes** panel shows the phenotype distribution as colors on the t-SNE plot. Finally, the **Feature ranking** panel shows a plot that ranks the discriminative power of feature sets (or gene sets) as the cumulative discriminant score for all phenotype variables.

The **Features** submodule performs clustering at either the gene level (**Gene** panel) or at the geneset level (**Geneset** panel). For both a gene or geneset UMAP plot is displayed, next to a "Gene Signatures" UMAP plot, where users can visualise specific phenotypes. Finally, below the plots, a table contains eithe the genes or genesets in a given selected area.

## **15.1 Samples**

### **15.1.1 Settings panel**

The settings panel on the right displays various options to customise the plots. Under `Show phenotypes` users can choose which phenotypes will be displayed in the Phenotype distribution plot under the **PCA/tSNE** panel. The `Split samples by` option applies only to heatmaps. "None" is the default view, "phenotypes" will redraw the heatmap based on the selected phenotype and "gene" allows users to split the heatmap based on the expression level of a gene that can be selected from a scrolldown menu. `Filter samples` can be used to represent only a specific subset of samples in the heatmaps, PCA/tSNE/UMAP plots and the parallel coordinates plot. Under `Gene family` a user can select whether the heatmap and parallel coordinates plot will be built based on all the genes, a specific contrast, a specific gene family or alternatively a custom list of genes provided by the user that can be pasted in the appropriate space. Furthermore, under *Advanced options*, users can choose the layout of the clustering plots (PCA, tSNE or UMAP), the level of analysis (gene or geneset) and exclude mitochondrial and ribosomal genes and/or genes in the X and Y chromosomes.

## 15.1.2 Heatmap

In the **Heatmap** panel hierarchical clustering can be performed on gene level or gene set level expression. For the latter, for each gene set (or pathway), a single-sample enrichment value is computed from the gene expression data using summary methods such as GSVA and ssGSEA.
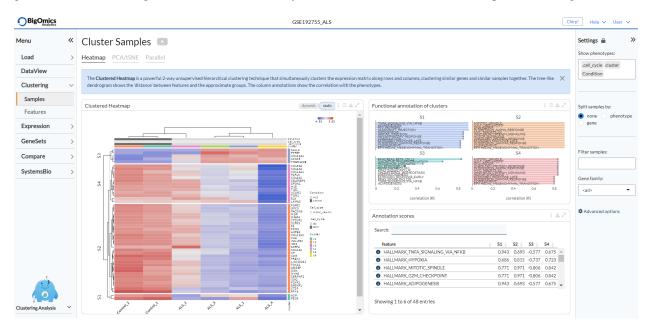
Next to the plot configuration settings, users can select between a "dynamic" or "static" heatmap. From the plot configuration settings on top of the plot, users can choose various options to customise their heatmaps. It is possible to order the top features under `top mode` as follows:

- sd - features with the highest standard deviation across all the samples,
- PCA - by principal components.
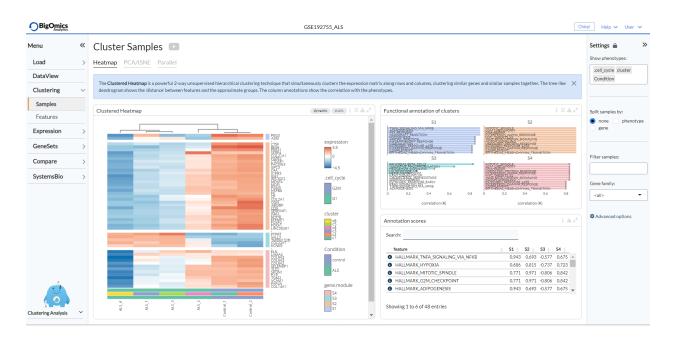- marker - features that are overexpressed in each phenotype class compared to the rest

In addition, users can specify the `Top N` (50, 150, 500) genes to be used in the heatmap and the number of gene clusters to be displayed under `K`. Users can also choose between 'relative', 'absolute' or 'BMC' (batch-mean centered) expression scale. Under the `CexCol` and `CexRow` settings, it is also possible to adjust the font sizes for the column and row labels. The legend in the heatmap can be disabled by unticking the `show legend` option.



The complex heatmap below is generated with the "static" option active. It is a clustered heatmap showing gene expression sorted by 2-way hierarchical clustering. Red corresponds to overexpression, blue to underexpression of the gene. At the same time, gene clusters are functionally annotated in the **Annotate clusters** panel on the right.



Activating the "dynamic" option generates an interactive version of the clustered heatmap. Users should be aware that for large datasets (such as single-cell RNA-seq data) this plot can become rather slow.

### 15.1.3 Annotate clusters

The features in the heatmap are divided into clusters depending on the selected `top mode` in the heatmap panel settings. For each cluster, the **Annotate cluster** section provides a functional annotation using more than 42 published reference databases, including but not limited to well-known databases such as MSigDB, Wikipathways, and GO. In the plot settings, users can specify the level and reference set to be used under the `Reference level` and `Reference set` settings, respectively. Users can also select a Fisher test weighting for gene sets.



The functional annotation for the clusters are displayed below, with the highest ranking annotation features (by correlation) displayed for each gene cluster. Length of the bars corresponds to the average correlation of the cluster with the annotation term. In the table below the barplots, users can view the correlation values of annotation features for each cluster.
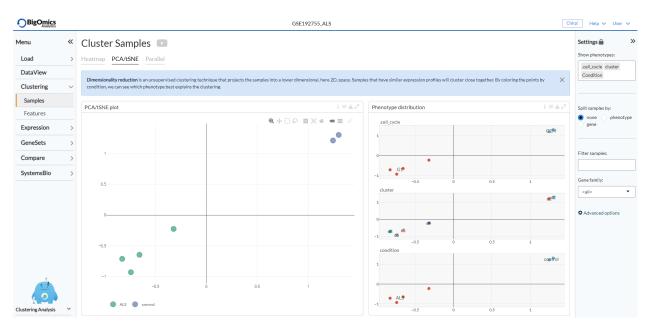
### 15.1.4 PCA/tSNE

The **PCA/tSNE** panel visualizes unsupervised clustering obtained by the principal components analysis (PCA), t-distributed stochastic embedding (tSNE) or Uniform Manifold Approximation and Projection (`UMAP <https://arxiv.org/abs/1802.03426>__`) algorithms. This plot shows the relationship (or similarity) between the samples for visual analytics, where similarity is visualized as proximity of the points. Samples that are 'similar' will be placed close to each other.

Users can customise the PCA/tSNE/UMAP plot in the plot settings, including the `color/label` and `shape` of points using a phenotype class, the placement of the plot legend at the bottom or as a group label, the inclusion of sample labels, the choice between a 2D/3D plot display and the normalisation of the plot matrix.

To the right of the PCA/tSNE/UMAP plot (labelled as *PCA/tSNE plot*) the platform also generates a group of plots, entitled *Phenotype distribution*, that visualise the distribution of the available phenotype data. The plots show the distribution of the phenotypes superposed on the t-SNE clustering.
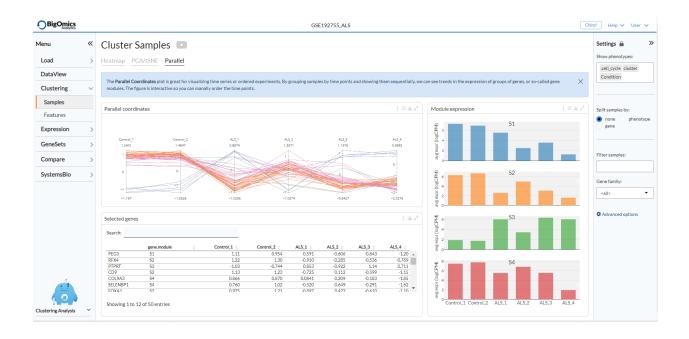


## 15.1.5 Parallel

The **Parallel** panel visualizes the expression levels of selected genes across all conditions in the plot labelled *Parallel coordinates*. The expression values are scaled but scaling can be removed via the plot settings, where gene expression levels can also be averaged by gene module. This interactive plot is particularly useful to users working with time series experiments, as samples can be grouped by condition (i.e. time) and ordered manually. A table (named *Selected genes*) containing average expression levels of selected genes across conditions is generated below the plot. Finally, to the right of the *Parallel coordinates* plot, a series of histograms are group together in the *Module expression* panel and display the overal expression of each module (the number of which is defined by the K value selected in the heatmap settings) by individual sample.
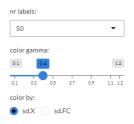
# 15.2 Features

## 15.2.1 Settings panel

Using the `Show phenotype` option the phenotypes that will be shown in the *Gene signatures* plots can be selected. Users can also select whether to use all genes for the *Gene UMAP* plot or instead select a combination of gene families under the `Annotate genes` option. The `Annotate genesets` option provides the same functionality for the *Geneset UMAP* plot based on the available geneset collections The `Show full table` option shows the full list of unfiltered genes or genesets. Under the advanced options users can select which sample group to use as a `Reference` to calculate the standard-deviation of log-expression (sd.X), or standard-deviation of the fold-change (sd.FC). If none is selected, the average values for all samples will be used instead. The final option, `UMAP datatype`, is used to select how the UMAP plot will be computed: either using the normalised log-expression (logCPM) or log-fold change matrix (logFC). logCPM is the default choice, while logFC can be used if batch or tissue effects are present in the dataset,
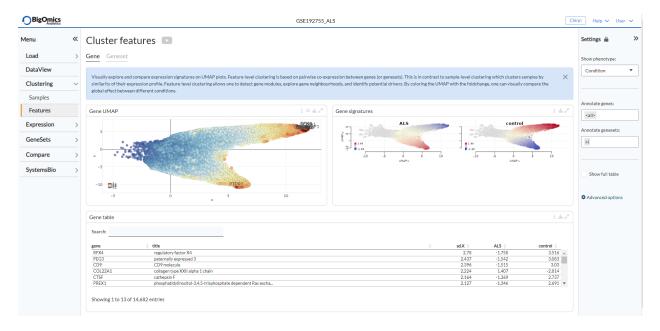
## 15.2.2 Gene and Geneset UMAP settings

Both the **Gene UMAP** and **Geneset UMAP** panels display the same setting options, as shown below. In it, users can use the `color by` option to select between standard-deviation of log-expression(sd.X) or standard-deviation of the fold-change (sd.FC) to colour individual genes or genesets. Users can also select the range of the colour intensity threshold using the `color gamma` option. Finally, users can choose the number of gene or genesets to be labelled in the UMAP plot under `nr labels`.
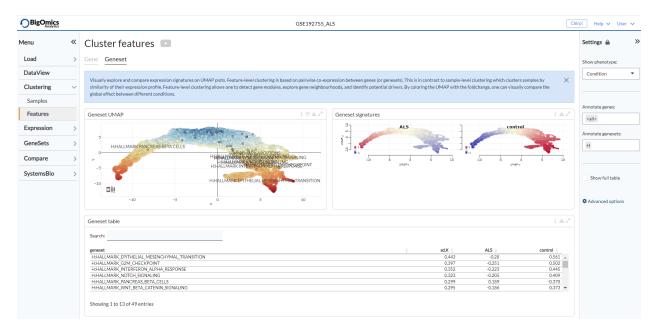


## 15.2.3 Gene

The **Gene** submodule contains three panels. The **Gene UMAP** panel displays the UMAP clustering of genes coloured by standard-deviation of log-expression(sd.X), or standard-deviation of the fold-change (sd.FC) and using the covariance of gene expression as a distance metric. The **Gene signatures** panel is to the right of the Gene UMAP, which shows the gene UMAP plots for each phenotypic group coloured by relative log-expression. Red indicates a gene is overexpressed in a specific group, blue that it is downregulated compared to the average values of all samples. The **Gene table** is located below the Gene UMAP and shows all the gene contained in the Gene UMAP plot and their relative change in expression. The contents of this table can be subsetted by selecting (by click&drag) on the Gene UMAP plot.

## 15.2.4 Geneset

The **Geneset** submodule contains three panels. The **Geneset UMAP** panel displays the UMAP clustering of genesets coloured by standard-deviation of log-expression(sd.X), or standard-deviation of the fold-change (sd.FC) and using the covariance of gene expression as a distance metric. The **Geneset signatures** panel is to the right of the Geneset UMAP, which shows the geneset UMAP plots for each phenotypic group coloured by relative log-expression. Red indicates a geneset is overexpressed in a specific group, blue that it is downregulated compared to the average values of all samples. The **Geneset table** is located below the Geneset UMAP and shows all the gene contained in the Geneset UMAP plot and their relative change in expression. The contents of this table can be subsetted by selecting (by click&drag) on the Geneset UMAP plot.

# EXPRESSION

The **Expression** module is divided into three submodules. The **Differential Expression** panel compares expression between two conditions (i.e. tumor versus control), which is one of the fundamental analysis in the transcriptomics data analytics workflow. For each comparison of two conditions (also called a 'contrast'), the analysis identifies which genes are significantly downregulated or overexpressed in one of the groups.

The **Correlation Analysis** submodule computes the correlation between genes and finds coregulated modules.

Finally, the **Find biomarkers** submodule is used to identify potential biomarkers for a set of phenotypic groups based on gene or protein expression values.

## 16.1 Differential expression

### 16.1.1 Settings panel

The settings panel on the right contains some settings for the analysis. Users can start the differntial expression (DE) analysis by selecting a contrast of their interest in the `Contrast` setting and specifying a relevent gene family in the `Gene family` setting. It is possible to set the false discovery rate (FDR) and the logarithmic fold change (logFC) thresholds under the `FDR` and `logFC threshold` settings, respectively. Under "Options", users can select to display all genes in the table (not only significant genes) and also select different combinations of statistical methods for the analysis.
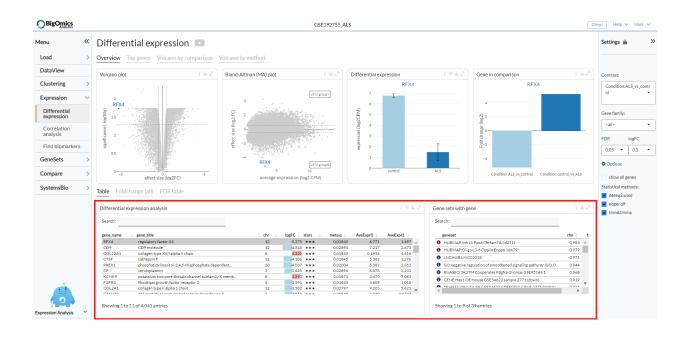
### 16.1.2 Table

The `Differential expression analysis` table shows the results of the statistical tests slected in the `Statistical methods`. By default, this table reports the meta (combined) results of DESeq2 (Wald), edgeR (QLF), and limma (trend) only. Users can filter top N = {10} differently expressed genes in the table by clicking the `top 10 genes` and also show the q-values from individual statistical methods by ticking `show individual q-values` from the table *Settings*.

For a selected comparison under the `Contrast` setting, the results of the selected methods are combined and reported in the `Differential expression analysis` table, where `meta.q` for a gene represents the highest `q` value among the methods and the number of stars for a gene indicate how many methods identified significant q values ($q < 0.05$). The table is interactive (scrollable, clickable); users can sort genes by `logFC`, `meta.q`, or average expression in either conditions.

By clicking on a gene in the the `Differential expression analysis` table (highlighted in grey in the figure), it is possible to see the correlation and enrichment value of gene sets that contain the gene in the the `Gene sets with gene` table.
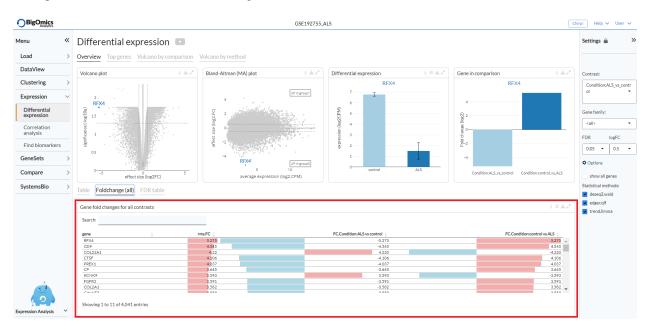
## 16.1.3 Foldchange (all)

The **Foldchange (all)** tab reports the gene fold changes for all contrasts in the selected dataset. The column `fc.var` corresponds to the variance of the fold-change across all contrasts.



## 16.1.4 FDR table

The **FDR table** tab reports the number of significant genes at different FDR thresholds for all contrasts and methods within the dataset. This enables to quickly see which methods are more sensitive. The left part of the table (in blue) correspond to the number of significant down-regulated genes, the right part (in red) correspond to the number of significant overexpressed genes.

## 16.1.5 Overview Plots

The **Overview** tab shows on top the following plots (from left to right):

**Volcano Plot**
> An interactive volcano plot for the chosen contrast. Unless a specific gene is selected from the differential expression analysis table, all significant genes are highlighted in blue.

**Bland-Altman (MA) plot**
> An interactive MA plot for the chosen contrast. Unless a specific gene is selected from the differential expression analysis table, all significant genes are highlighted in blue.

**Differential expression**
> Differential expression boxplot for a gene that is selected from the differential expression analysis table. Users can customise the plot via the settings icon on top to ungroup samples and change the scale to counts per million (CPM).

**Gene in comparison**
> Fold change summary barplot across all contrasts for a gene that is selected from the differential expression analysis table.



## 16.1.6 Top genes

The **Top genes** tab shows the average expression plots across the samples for the top differentially (both positively and negatively) expressed genes for the selected comparison from the `Contrast` setting.

The plot can be customised via the settings to remove the log scale, group samples and show only samples included in the selected contrast.

### 16.1.7 Volcano by comparison

Under the **Volcano by comparison** tab, the platform simultaneously displays multiple volcano plots for genes across all contrasts. By comparing multiple volcano plots, the user can immediately see which comparison is statistically weak or strong. Experimental contrasts with better statistical significance will show volcano plots with 'higher' wings.
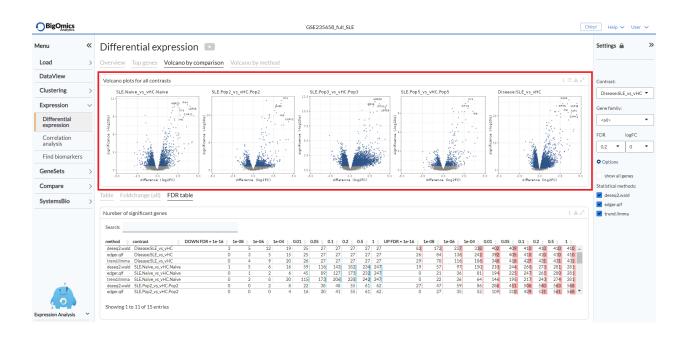
### 16.1.8 Volcano by method

Under the **Volcano by method** tab, the platform displays the volcano plots provided by multiple differential expression calculation methods for the selected contrast. Methods showing better statistical significance will show volcano plots with 'higher' wings.

## 16.2 Correlation analysis

### 16.2.1 Settings panel

The panel contains the main settings for the analysis. The analysis can be started by selecting a gene of interest from the `Gene` settings. Users can filter for a specific gene family by using the `Filter genes` setting. Finally, users can select the number of top genes used to compute partial correlation under the homonymous setting.

## 16.2.2 Correllation

Under the **Correlation** tab, the platform displays three different outputs:

**Top correlated genes**
   A plot displaying the highest correlated genes in respect to the selected gene.

**Correlation scatter plots**
   Pairwise scatter plots for the co-expression of correlated gene pairs across the samples.

The straight line correspond to the (linear) regression fit. Using the settings on top, the plot can be customised by changing the colour of the gene pairs by phenotype using the `Colour by` option. Users can also change the layout of the plots under ``Layout`.



**Correlation table**
   it shows the statistical results from correlated gene pairs.

## 16.2.3 Graph

Two plots are generated under the **Graph** tab:

**Partial correlation network**
   This plot shows the correlations between gene pairs. Grey edges correspond to positive correlation, red edges correspond to negative correlation. The width of the edge is proportional to the absolute partial correlation value of the gene pair. It can be customised from the plot settings, where the radius of the plot and the partial correlation threshold (`pcor threshold`) can be altered. Users can also the choose from four different layouts (Fruchterman-Reingold, Kamada-Kawai, graphopt and tree layout) under the `Layout` option.

**Correlation UMAP**

UMAP plot showing the correlation between genes. Genes that are correlated are generally positioned close to each other. Red corresponds to positive correlation/covariance, blue to negative. User can select whether to colour the genes by correlation or covariance in plot settings.



The **Partial correlation network** is on the left, while the **Correlation UMAP** is on the right of the page.



# 16.3 Find Biomarkers

## 16.3.1 Settings panel

The panel contains the main settings for the analysis. Users select the phenotype on which to perform the analysis with the `Predicted target` setting. In the `Filter samples` setting users can restrict the analysis to a subset of samples rather than the default choice of all samples. The `Feature set` setting will by default be set to "all", but users can restrict the calculations only to a specific gene family or, alteratively, they can paste their own selection of genes. Once all the settings have been assigned, clicking *Compute* will start the calculation.

## 16.3.2 Feature Selection

Four panels (from left to right and top to bottom) are present under the **Feature Selection** tab:

**Variable Importance**

Omics Playground calculates a variable importance score for each feature using multiple state-of-the-art machine learning algorithms, including LASSO, elastic nets, random forests, and extreme gradient boosting. Note that we do not use the machine learning algorithms for prediction but we use them just to compute the variable importances according to the different methods. An aggregated score is then calculated as the cumulative rank of the variable importances of the different algorithms.

By combining several methods, the platform aims to select the best possible predictive features. The top features are determined as the features with the highest cumulative ranks. .

**Biomarker Expression**
These boxplots shows the expression of putative biomarkers across the samples of the identified features.

**Heatmap**
Expression heatmap of top gene features according to their variable importance. By default the plot will only show the selected samples, but all samples can be displayed via the settings button on top.

**Decision tree**
The decision tree shows a tree solution for classification based on the top most important features. The plot provides a proportion of the samples that are defined by each biomarker in the boxes.

### 16.3.3 Feature-set Ranking

This tab consists of a single plot, the aptly named **Feature-set Ranking**, which displays the ranked discriminant score for top feature sets. The plot ranks the discriminative power of the feature set (or gene family) as a cumulative discriminant score for all phenotype variables. In this way, we can find which feature set (or gene family) can explain the variance in the data the best.

Users can choose between three different methods for the calculation of the plot. P-value based scoring ('p-value') is computed as the average negative log p-value from the ANOVA test. Correlation-based discriminative power ('correlation') is calculated as the average 1-cor between the groups. Thus, a feature set is highly discriminative if the between-group correlation is low. The 'meta' method combines the scores of the two aforementioned methods in a multiplicative manner.

# GENESETS

This module is subdivided into four submodules: **Geneset enrichment**, **Test Geneset**, **Pathway Analysis** and **Word cloud**.

## 17.1 Geneset enrichment

Similar to the differential gene expression analysis, users can perform differential expression analysis on a geneset level in this page, which is also referred as gene set enrichment (GSE) analysis. The platform has more than 50.000 genesets (or pathways) in total that are divided more than 30 geneset collections such Hallmark, Reactome, Wikipathways, and Gene Ontology (GO). Users have to specify which comparison they want to visually analyze employing a certain geneset collection.

The tab consists of the settings panels as well as four results panels: **Top enriched gene sets**, **Frequency in top gene sets**, a **Table** panel where users can swicth between three tables and **Genes in gene sets**.
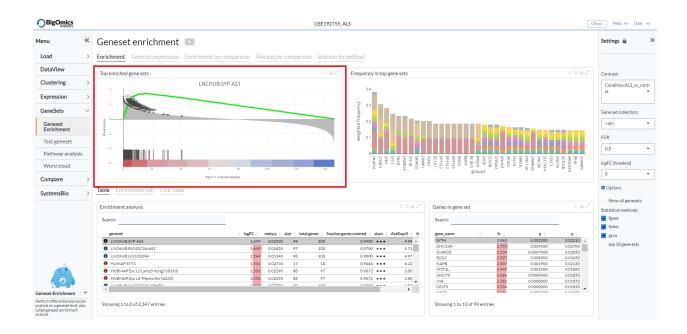
### 17.1.1 Settings panel

The enrichment analysis can be started by selecting a contrats of interest in the `Contrast` and specifying a relevent gene set family in the `Gene set collection`. It is possible to set the false discovery rate (FDR) and the logarithmic fold change (logFC) thresholds under the `FDR` and `logFC threshold` settings, respectively. We allow users to show all gene sets in the``Enrichment analysis`` table , rather than just the significant ones, under the *Options* menu. Users can also select statistical methods for the enrichment analysis from the same menu. To ensure the statistical reliability, the platform performs enrichment analyses using multiple methods, including Spearman rank correlation, GSVA, ssGSEA, Fisher's exact test, GSEA, camera and fry. Finally users can limit the results to only the top 10 up- and down-regulated gene sets.

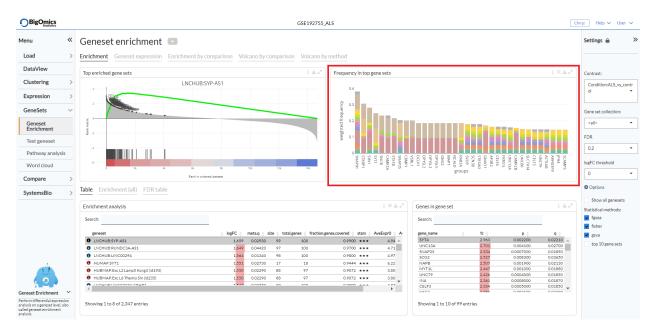### 17.1.2 Top enriched gene sets

This panel displays the enrichment plots of the top differentially enriched gene sets (up and down) for the selected contrast. Black vertical bars indicate the rank of genes in the gene set in the sorted list metric. The green curve corresponds to the enrichment score (ES). The more the green ES curve is shifted to the upper left of the graph, the more the gene set is enriched in the first group. Conversely, a shift of the ES curve to the lower right, corresponds to more enrichment in the second group. This panel will by default display the 12 most up-regulated genesets in the selected contrast. Selecting a specific gene set from the 'Enrichment analysis' table below it will display the selected gene set alone.

### 17.1.3 Frequency in top gene sets

This panel shows the number of times a gene is present in the top-N genesets sorted by frequency. Genes that are frequently shared among the top enriched gene sets may suggest driver genes.



The settings icon open the settings menu from which users can select the number of top gene sets to be used (`Number of top sets`), whether to weight by geneset size (`Weight by geneset size`) and whether to weight by fold-change (`Weight by FC`).
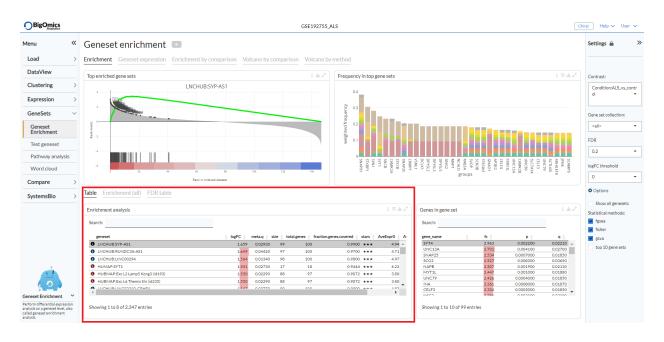


### 17.1.4 Table

This table shows the enrichment analysis results from the selected statistical methods. With default settings, this table reports the meta (combined) results of camera, fgsea, and Spearman rank correlation only. Users can also display individual q-values for each of the selected analysis methods using the settings on top of the table.
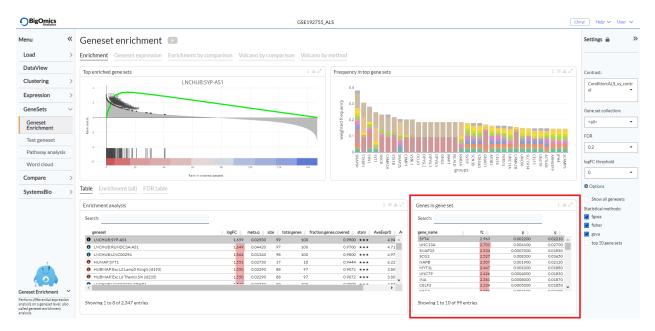


The combined enrichment analysis results from the methods are displayed in this table, where for each geneset the `meta.q` corresponds to the highest `q` value provided by the methods and the number of `stars` indicate how many methods identified the geneset as significant (`q < 0.05`). `AveExpr0` and `AveEprx1` refer to the average gene set expression in group 0 and group 1 of the selected pairwise comparison respectively. The table is interactive; users can sort it by `logFC`, `meta.q`, `AveExpr0`, `AveEprx1` and `stars`.

### 17.1.5 Genes in gene set

By clicking on a gene set in the **Table** above, it is possible to see the gene list of that gene set under **Genes in gene set**. this table also reports the fold-change, statistics and correlation of the genes in the selected gene set. By clicking on a gene in this table, users can check the expression status of the gene for the selected contrast in the *Expression barplot* and its correlation to the gene set in the *Enrichment vs. expression* scatter plot under the **Geneset expression** tab, discussed in the next section.

## 17.1.6 Geneset expression

The **Geneset expression** panel provide plots associated with the gene set selected in **Table** and gene selected in **Genes in gene set**, as explained in the previous section.

**Volcano Plot**

Volcano-plot of genes showing the significance versus the fold-change on the y and x axes, respectively. Genes in the selected gene set are highlighted in blue.
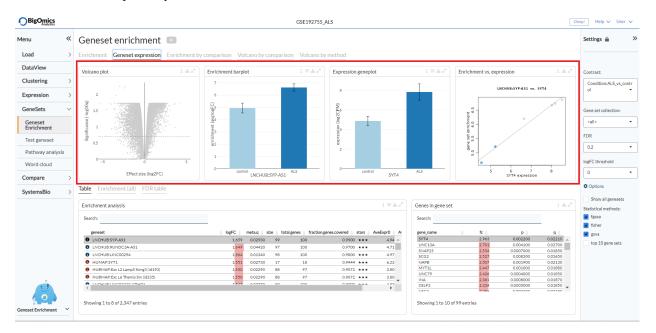
**Enrichment barplot**

Enrichment plot associated with the gene set selected from **Table**. Samples in the barplot can be ungrouped via the plot settings icon.

**Expression barplot**

Barplot of the gene expression of the gene. Samples in the barplot can be ungrouped via the plot settings icon.

**Enrichment vs. expression**

Scatter plot of the enrichment versus the expression of the selected geneset and gene, on the y and x axes, respectively.



## 17.1.7 Enrichment by comparison

Under the **Enrichment by comparison** panel, users can fin enrichment plots for the selected gene set (in **Table**) across multiple contrasts. The figure allows to quickly compare the enrichment of a certain gene set across all other comparisons.

## 17.1.8  Volcano by comparison

The **Volcano by comparison** panel simultaneously displays volcano plots of gene sets enrichment across all contrasts, showing the enrichment score versus significance on the x and y axes, respectively. This provides users an overview of the statistics across all comparisons. By comparing multiple volcano plots, the user can immediately see which comparison is statistically weak or strong. Experimental contrasts showing better statistical significance will show volcano plots with 'higher' wings.

## 17.1.9 Volcano by method

Under the **Volcano by method** panel, users can see the simultaneous visualisation of volcano plots of gene sets for different enrichment methods. This provides users an quick overview of the sensitivity of the statistical methods at once. Methods showing better statistical significance will show volcano plots with 'higher' wings.



## 17.1.10 Enrichment (all)

The **Enrichment (all)** table provides the enrichment scores of gene sets across all contrasts. It is also possible to visualise all q-values from the setting icon on top of the plot.
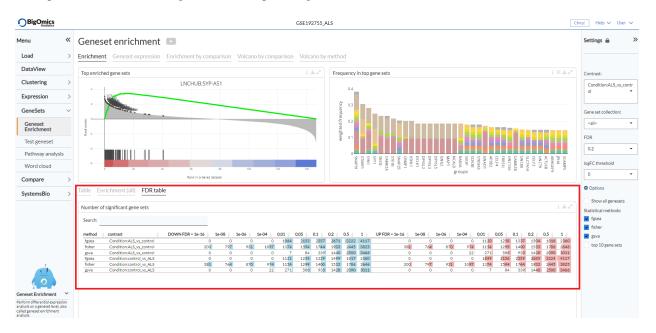
### 17.1.11 FDR table

The **FDR table** panel reports the number of significant gene sets at different FDR thresholds for all contrasts and methods. Using the table the user can determine which statistical methods perform better for a particular contrast. The left part of the table (in blue) correspond to the number of significant down-regulated gene sets, the right part (in red) correspond to the number of significant overexpressed gene sets.

## 17.2 Test geneset

The **Test geneset** submodule is used to test a specific user-defined geneset, or, alternatively, a geneset from the KEGG or Hallmark collection or a list of genes extracted from the available pairwise comparisons, for enrichment in the experimental dataset. The plot contain five tabs: **Enrichment table**, **Volcano plots**, **Enrichment**, **Overlap/similarity** and **Markers**.

### 17.2.1 Settings panel

The settings panel of the **Test geneset** module consists of a box where users can paste gene lists to be queried. Three examples of gene lists are provided below it, namely the *apoptosis*, *cell_cycle* and *immune_chkpt* lists. Users can also select specific KEGG and hallmark gene sets from the *Options* in the panel, as well as using the differentially expressed genes or proteins from one of the contrasts.

## 17.2.2 Enrichment table

This tab is shown indepently of the other four. It display two tables. The first table, *Enrichment by contrasts*, shows the enrichment scores of query signature across all contrasts. The table summarizes the enrichment statistics of the gene list in all contrasts using the GSEA algorithm. The NES corresponds to the normalized enrichment score of the GSEA analysis.
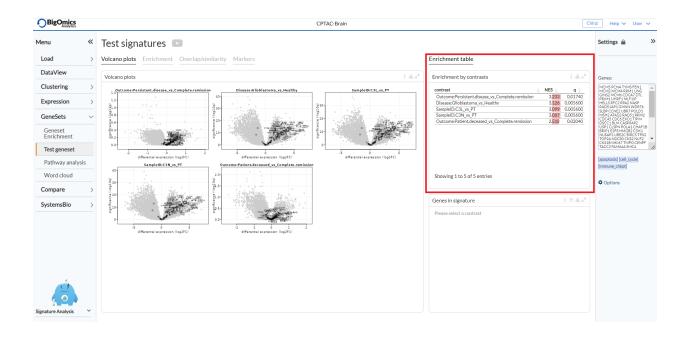
The second table, *Genes in signature*, is based on which pairiwse comparison is selected in *Enrichment by contrasts* and shows the genes of the current signature corresponding to the selected contrast. Genes are sorted by decreasing (absolute) fold-change.
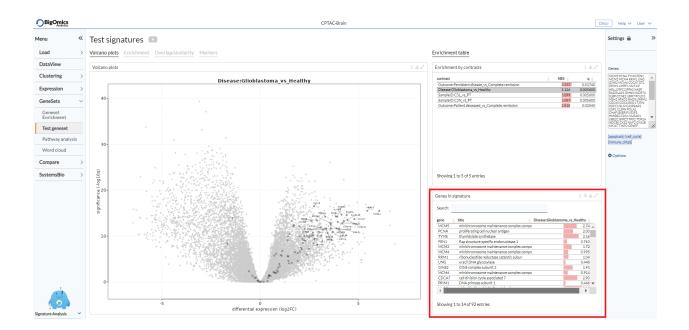
## 17.2.3 Volcano plots

This tab contains volcano plots showing where the genes of a test signatures fall in the experimental pairiwse comparisons. For positive enrichment, genes of the query signature would fall on the upper right of the volcano plot, for negative enrichment, on the upper left.

## 17.2.4 Enrichment

In this tab, multiple plots show the enrichment of the query signature in all constrasts. Positive enrichment means that a particular contrast shows similar expression changes as the query signature. Users can either view the plots of all teh contrasts, while selecting a specific contrast in the *Enrichment by contrasts* table will only display the ernichment plot for it.

## 17.2.5 Overlap/similarity

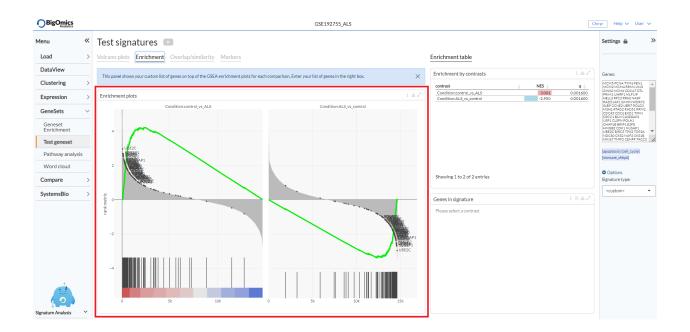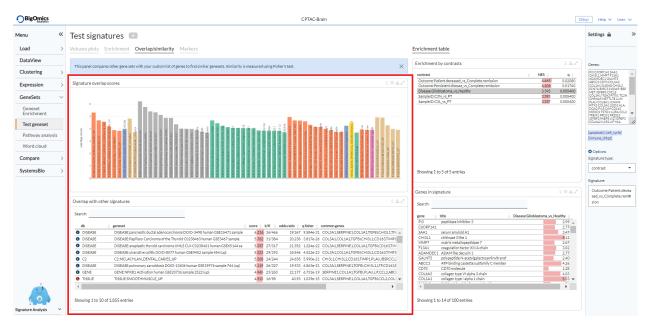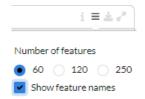Under the Overlap/similarity tab, users can find the similarity of their gene list with all the gene sets and pathways in the platform, including statistics such as the total number of genes in the gene set (K), the number of intersecting genes between the list and the gene set (k), the overlapping ratio of k/K, logarithm of the odds ratio (log.OR), as well as the p and q values by the Fisher's test for the overlap test.

The *Signature overlap scores* plot is located on top. Its vertical axis shows the overlap score of the gene set which combines the odds ratio and significance (q-value) of the Fisher's test. The*Overlap with other signatures* table is located below it and contain various statistical values.



The *Signature overlap scores* plot can also be customised via the settings icon, where users can select the number of signatures to display (`Number of features`) and can toggle the signature names on and off (`Show features names`).



## 17.2.6 Markers

The **Markers** tab contains a series of plots (under *Markers plot*) showing the expression levels of the tested genes in the dataset samples as a colored t-SNE plot in red (highly expressed) and light grey (low expressed). The first figure shows the single-sample enrichment of your signature list in red (upregulation) and blue (downregulation).

The plots can be sorted via the `Sort by` option in the settings icon by correlation (default), probability or name. `Layout` swaps the layout of the plots between a 4x4 (default) and a 6x6 grid.

## 17.3 Pathway analysis

This module performs specialized enrichment analysis providing higher level functional and visual interpretation.

The **WikiPathways** panel maps the differential fold-changes onto the WikiPathways pathway maps. The **Reactome** panel does the same for the Reactome pathway maps. Under the **GO** panel, a graph-based enrichment analysis is done using the Gene Ontology (GO) graph structure.

### 17.3.1 Settings panel

Users can specify the contrast of their interest in the `Contrast` settings. Under the main *Options*, users can select `filter significant (tables)` to keep only significant entries in the table.

### 17.3.2 WikiPathways

WikiPathways is a public collection of manually curated pathways representing the current knowledge of molecular interactions, reactions and relation networks as pathway maps. In the **WikiPathways** panel, each pathway is scored for the selected contrast profile and reported in the table. An activation-heatmap comparing the activation levels of pathways across multiple contrast profiles is generated. This facilitates the quick detect the similarities between profiles in certain pathways. More detailed explaination of each output is provided below.

> **WikiPathway**
>
> In the pathway map, genes are colored according to their upregulation (red) or downregulation (blue) in the contrast profile. Each pathway is scored for the selected contrast profile and reported in the table below the plot.
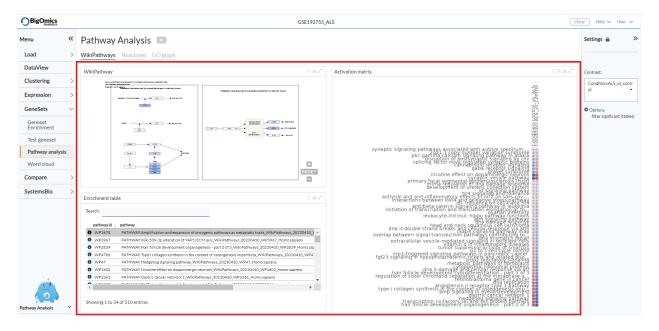
### 17.3.3 Reactome

Reactome is a collection of manually curated pathways representing the current knowledge of molecular interactions, reactions and relation networks as pathway maps. In the **Reactome** panel, each pathway is scored for the selected contrast profile and reported in the table. A unique feature of the platform is that it provides an activation-heatmap comparing the activation levels of pathways across multiple contrast profiles. An activation-heatmap comparing the activation levels of pathways across multiple contrast profiles is generated. More detailed explaination of each output is provided below.
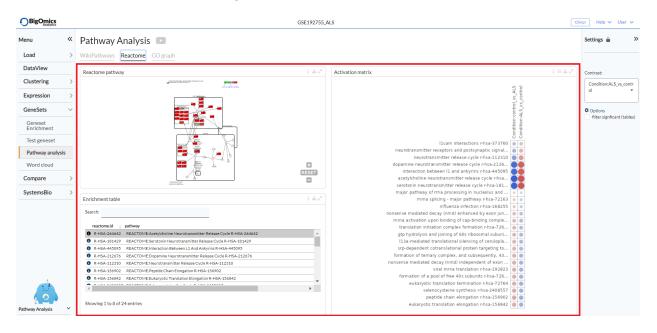
**Reactome**
> In the pathway map, genes are colored according to their upregulation (red) or downregulation (green) in the contrast profile. Each pathway is scored for the selected contrast profile and reported in the table below the plot.

**Enrichment table**
> The table is interactive; enabling user to sort on different variables (Reactome id, pathway, logFC and meta q-values) and select a pathway by clicking on the row in the table.

**Activation matrix**
> The activation matrix visualizes the activation levels of pathways (or pathway keywords) across multiple contrast profiles. This facilitates to quickly see and detect the similarities of certain pathways between contrasts. The size of the circles correspond to their relative activation, and are colored according to their upregulation (red) or downregulation (blue) in the contrast profile. The matrix can be normalised from the *settings* icon.
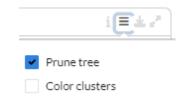
## 17.3.4 GO graph

In the **GO** panel, users can perform GO analysis. GO defines functional concepts/classes and their relationships as a hierarchical graph. The GO database provides a computational representation of the current knowledge about roles of genes for many organisms in terms of molecular functions, cellular components and biological processes. All the features described under the **KEGG pathway** panel, such as scoring the gene sets and drawing an activation-heatmap, can be performed for the GO database under the GO graph tab. Instead of pathway maps, an annotated graph structure provided by the GO database is potted for every selected gene set. Each output chart/table of the panel is describer below in detail.

**Gene Ontology graph**

The structure of GO can be described in terms of a graph, where each GO term is a node, and the relationships between the terms are edges between the nodes. GO is loosely hierarchical, with 'child' terms being more specialized than their 'parent' terms. The graph is interactive. You can move the graph and zoom in using the mouse. Under the graph *settings*, users can select `Prune tree` to prune the tree only with significant branches and `color custers` to highlight clusters with different colors



**GO score table**

The scoring of a GO term is performed by considering the cumulative score of all terms from that term to the root node. That means that GO terms that are supported by higher level terms levels are preferentially scored.

**Activation matrix**

The GO activation matrix visualizes the activation of GO terms across conditions. From this figure, you can easily detect GO terms that are consistently up/down across conditions. The size of the circles correspond to their relative activation, and are colored according to their upregulation (red) or downregulation (blue) in the contrast profile. The matrix can be normalised from the *settings* icon.

## 17.4 Word Cloud

The **WordCloud** panel performs "keyword enrichment analysis". It computes enrichment of a selected keyword across all contrasts. Select a keyword by clicking a word in the 'Enrichment table'. Keyword enrichment is computed by running GSEA on the enrichment score profile for all contrasts. We defined the test set as the collection of genesets that contain the keyword in the title/description.

### 17.4.1 Settings panel

Users can specify the contrast of their interest in the `Contrast` settings.



### 17.4.2 Main panel

The main panel consists of five different outputs:

**Enrichment plots**
> The plots visualize the enrichment of the selected keyword in the contrasts. Black vertical bars indicate the position of gene sets that contains the *keyword* in the ranked list of enrichment scores. The curve in green corresponds to the 'running statistic' of the keyword enrichment score. The more the green ES curve is shifted to the upper left of the graph, the more the keyword is enriched in the first group. Conversely, a shift of the green ES curve to the lower right, corresponds to keyword enrichment in the second group.

**Word cloud**
> The word cloud plot visualizes the frequency/enrichment of keywords for the data set. Select a keyword in the 'Enrichment table'. In the plot settings, users can exclude certain words from the figure, or choose the color palette. The sizes of the words are relative to the normalized enrichment score (NES) from the GSEA computation. Keyword enrichment is computed by running GSEA on the mean (squared) enrichment profile (averaged over all contrasts). For each keyword, we defined the 'keyword set' as the collection of genesets that contain that keyword in the title/description.

**Word t-SNE**

The Word t-SNE plot visualizes the similarity of the keywords that were found in the title/description of gene sets. Keywords that are often found together in title/descriptions are placed close together in the t-SNE. For each keyword we computed enrichment using GSEA on the mean (absolute) enrichment profiles (averaged over all contrasts). Statistically significant gene sets (q<0.05) are colored in red. The sizes of the nodes are proportional to the normalized enrichment score (NES) of the keyword. In the plot settings, the user can choose between t-SNE and "Uniform Manifold Approximation and Projection" (UMAP).



**Enrichment table**

The Enrichment table summarizes the results from the enrichment test for the tested keywords. The NES corresponds to the normalized enrichment score from the GSEA analysis.

**Leading-edge table**

The Leading-edge table shows the geneset titles that have contributed to the enrichment of the selected keyword.

# COMPARE

The **Compare** module consists of three submodules: **Compare signatures**, **Compare datasets** and **Similar experiments**.

**Compare signatures** allows users to compare experiments by intersecting their signature genes. Under **Compare datasets**, users can compare pairwise comparisons from the dataset in use and another dataset stored in the platform. Finally, with the **Similar experiments** panel, users can perform a large-scale comparison across all uploaded datasets or a pre-loaded collection of more than 6000 public dataset.

## 18.1 Compare signatures

The main goal of this submodule is to identify contrasts showing similar profiles and find genes that are commonly up/down regulated between two contrasts. The panel is divided into two subpanels: **Pairwise scatter** and **Signature clustering**.

### 18.1.1 Settings panel

Users can select contrasts to compare from the `Contrast` settings in the input panel on the left. Under *Options*, the `Level` setting allows users to toggle between gene or gene set level analysis, while users can set the `Filter` for selecting specific features (e.g. a specific gene family or gene set).

### 18.1.2 Pairwise scatter

The **Pairwise scatter** panel shows three outputs:

**Scatterplot pairs**
    The Pairs plot provides interactive pairwise scatterplots for the differential expression profiles corresponding to multiple contrasts. The main purpose of this panel is to identify similarity or dissimilarity between selected contrasts. When K>=3 contrasts are selected, the figure shows a KxK scatterplot matrix. Via the *Settings*, users can disable the highlighting of genes on the plot.

**Venn diagram**
    The Venn diagram visualizes the number of intersecting genes between the selected contrast profiles. The diagram can be customised via the settings icon by altering `logFC` and `FDR` thresholds and by selecting whether to view disregulated genes jointly or separated by up- and down-regulation under `Counting`.

**Leading-edge table**
    Venn diagram areas can be selected via the settings icon (`Filter intersection`) and are represented by corresponding letters (e.g. 'ABC' represents the intersection of contrasts A, B and C). Contrast letter identifiers are provided in the Venn Diagram.

Settings 🔒                    »

Contrasts:

Disease:Glioblastoma_v
s_Healthy

Outcome:Patient.decea
sed_vs_Complete.remis
sion

Outcome:Persistent.dis
ease_vs_Complete.remi
ssion

⚙ Options

Level:

● gene      ○ geneset

Filter:

<all>              ▼

FDR          logFC

0.5    ▼    0.2    ▼

Counting:

● both    ○ up/down

Filter intersection:

ABC                ▼

The three output panels are highlighted in the figure below.



### 18.1.3  Signature clustering

Two plots are showed in this panel:

**Foldchange heatmap**
The foldchange heatmap shows the similarity of the contrasts visualized as a clustered heatmap. Contrasts that are similar will be clustered close together. The plot can be customised via the settings icon. Users can select to show only the selected contrasts (default is to show all), cluster the genes on the heatmap and change the annotation type between logFC (*boxplot*) and cumulative logFC (*barplot*).



**Contrast correlation**
Contrasts that are similar will be clustered close together. The numeric value in the cells correspond to the Pearson correlation coefficient between contrast signatures. Red corresponds to positive correlation and blue to negative correlation. Under the plot settings, users can use `show all contrasts` (default) to show all contrasts or only the selected ones and change the `number of top genes` to specify the number of top genes for computations (default=1000).

A typical output can be seen below.



## 18.2 Compare datasets

With this submodule, users can compare pairwise comparisons across datasets that have been uploaded into the platform. The submodule is split into thee tabs: **Compare expression**, **Foldchange** and **Gene Correlation**.

### 18.2.1 Settings panel

Users can select the pairwise comparisons to be selcted from `Dataset1` and `Dataset2`, from which they can also select a dataset from the list of uploaded experiments. Under *Options*, users can set the `Plot type` (default: UMAP1) to be displayed in the **Dataset1** and **Dataset2** panels. `Highlight genes` is used to label genes in the **Dataset1** and **Dataset2** plots. The choice is between highlighting top genes/proteins (default) or a custom list that users can type or paste in the corresponding box. With `ntop` users can define how many genes or proteins to label.

### 18.2.2 Compare expression

This panel shows plots for selected pairwise comparisons from the current dataset (*Dataset1*) and a second dataset (*Dataset2*) selected from a list of uploaded experiments. The type of plot can be selected via the main submodule **Settings**. Users can select between plotting the genes or proteins as UMAP plots based on either dataset 1 or dataset 2 (UMAP1 and UMAP2), as volcano, MA plots, scatter plots or heatmaps.

## 18.2.3 Fold change

This tab contains three panels: **FC correlation**, which contains the plot between the pairwise comparsions selected between two datasets and the **Cumulative foldchange** panels that show barplot highlighting the fold changes in expression for each of the selected pairwise comparisons.

> **FC correlation**
> Scatter plot of gene expression scatter values between two contrasts. Scatters that are similar show high correlation, i.e. are close to the diagonal.

> **Cumulative foldchange upper**
> Barplot showing the cumulative fold changes on dataset 1.

> **Cumulative foldchange lower**
> Barplot showing the cumulative fold changes on dataset 2.



## 18.2.4 Gene correlation

The **Gene Correlation** tab is used to compare the expression levels of individual genes or proteins between pairwise comparions across datasets. It can also be used for the combined analysis of proteomics and transcriptomics datasets. It contains three panels: **Expression**, **Correlation score** and **Gene correlation**.

> **Expression**
> Barplots of expression values for multiple comparisons in the two datasets (blue and green). Bars are labelled by pairwise comparison groups.

> **Correlation score**
> In this searchable table, users can check mean expression values and correlation scores of genes/proteins across the selected pairwise comparisons.

> **Gene correlation**
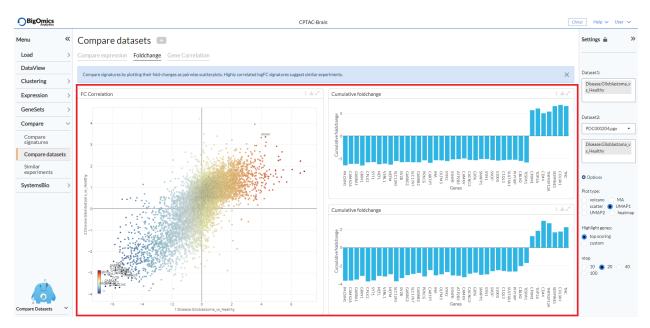> Scatter plots of gene expression scatter values between two contrasts. Scatters that are similar show high correlation, i.e. are close to the diagonal. This plot is only available for studies with matched sample Ids and can be used to compare proteomics and transcriptomics datasets from the same samples. Users can select by which pairwise comparison to colour the samples via the settings icon.

The genes or proteins appearing in the **Expression** barplots and **Gene correlation** scatter plots will be determined in the **Correlation score** table.



## 18.3 Similar experiments

With the final submodule, users can compare contrasts from different experiments against a selected pairwise comparisons from the current dataset simultenously. Rather than a pairwise analysis between two datasets, this analysis will take into accounts all uploaded datasets or, alternatively, access two databases of more than 6000 individual experiments collected from the GEO database. The submodule is split into three tabs: **FC correlation**, **FC heatmap** and **Meta-network**.

### 18.3.1 Settings panel

Through the **Settings** panel, an experimental `Contrast` from the selection of available pairwise comparisons in the dataset can be set. Under `Signature` DB, users can select from three databases: *datasets-sigdb.h5* corresponds to the datasets that have been uploaded by teh user into the platform, *sigdb-archs4.h5* correspond to a collection of more than 6000 datasets from the GEO database, while *sigdb-virome.h5* contains a selection of more than 300 viral infection GEO datasets. Under the *Advanced options*, users can choose whether to `hide cluster contrasts` (default) and show the absolute score (`abs.score`), which is the default choice. Under `Select genes` it is possible to type or paste a user-generated list of genes/proteins and also select the number of genes to be labelled (default: 50). Clicking `Recalculate` will then generate new plots based on the inputed gene list.

## 18.3.2  FC correlation

With this tab it is possible to compare different experiments by correlating their fold-change signatures. The tab consists of three panels: **FC scatter plots**, **Similarity scores** and **FC-FC scatterplot**.

**FC scatter plots**

Scatter plots of gene expression foldchange values between two contrasts. Foldchanges that are similar show high correlation, i.e. are close to the diagonal. You can switch to gsea or UMAP enrichment plots in the `plot type` option in the settings icon.



**Similarity score**

In this searchable table, Normalized enrichment scores (NES) and Pearson correlation (rho) of reference profiles with respect to the currently selected contrast are displayed. The top 100 up/down genes are considered for the calculation of rho or NES. The score is calculated as rho^2*NES. Highlighting a specific dataset will change the FC-FC scatterplot accordingly.

**FC-FC scatterplot**

This plot provides a pairwise scatterplot of logFC fold-change profiles for the selected contrasts. The main purpose of this panel is to identify similarity or dissimilarity between selected contrasts. The scatter plot is interactive and shows information of each gene by hovering over it with the mouse. The `logFC threshold` can be set via the settings icon.



The three panels are highlighted in the image below.

## 18.3.3 FC heatmap

The **FC heatmap** tab provides three panels for more comparative analysis: **Connectivity map**, **Fold-change table** and **Connectivity heatmap**.

**Connectvity map**

The connectivity map plot shows the similarity of logFC signatures as a t-SNE plot. Signatures that are similar will be clustered close together, signatures that are different are placed farther away. The plot can be extensively customised via the options in the settings icon. `Layout` is usde to alter the layout of the plot, namely pca, tsne (default) and volcano. `Score threshold` thresholds the points by minimum score.``Color by`` is used to colour the samples by either score (default) or dataset. `Color gamma` is used for colour adjustment. Under `Other options`, users can add lables to the plot, turn it into a 3D plot, change the background colour to black and increase the size of the dots.

Plot type:

scatter  gsea  umap

**Fold-change table**

In this searchable table, gene expression fold-changes (log2FC) of similar signatures across different experiments are indicated.

**Connectivity heatmap**

The heatmap clusters contrasts that are similar close together to provide an overview of the top most correlated contrasts with th equeried pairwise comparison. Through the settings icon, users can alter the `Number of signatures` to be displayed (default=20), they can also `Cluster genes` rather than sorting them by expression (default: off), `Use absolute foldchange` for calculating the cumulative sum and finally can `Reverse negative contrasts` (default: on), so as not to consider the direction of the correlation but only its strength.

These plots are complementary to the plots provided by the **FC correlation** tab, to explore more in detail the nature of the correlation between contrasts.

## 18.3.4 Meta-network

The final tab of the **Similar experiments** submodule provides some comparative gene network analysis as well as geneset level annotation of the most frequently enriched genesets across the correlated pairwise comparisons. It consists of four panels: **Leading-edge graph**, **Cumulative foldchange**, **Enrichment graph** and **Cumulative enrichment**.

**Leading-edge graph**

This graph displays the connections between genes shared across the correlated signatures. The edge width corresponds to the number of signatures that share that pair of genes in their top differentially expressed genes. Through the settings icon, it is possible to set the `Edge threshold`, select the number of signatures (`nr of signatures`) to be used (default=10) and select the `Size` parameter for the nodes (default: cumFC).

**Cumulative foldchange**

The barplot visualizes the cumulative foldchange between the top-10 most similar profiles. Genes that are frequently shared with high foldchange will show a higher cumulative score. You can choose between signed or `Absolute foldchange` in the options and use `order by` to order the plot by FC or cumulative FC (cumFC, default).

**Enrichment graph**

In this graph, the edge width corresponds to the number of signatures that share that pair of genesets in

Layout:

○ pca   ● tsne   ○ volcano

Score threshold:

[0]                                    [1]

0  0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

Color by:

● score   ○ dataset

Color gamma:

[0.1]  [0.5]                          [2]

0.1 0.3 0.5 0.7 0.9 1.1 1.3 1.5 1.7 1.9 2

Other options:

☐ show label

☐ 3D plot

☐ dark mode

☐ larger points

edge threshold:

| 0 | 0.2 | | 1 |

0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1

nr of signatures:

○ 5  ● 10  ○ 25  ○ 100

Size:

○ FC  ● cumFC  ○ centrality

☐ Absolute foldchange

Order:

○ this FC  ● cumFC

their top enriched genesets. In the plot options you can set the `Edge threshold`, the number of similar experiments to consider (`N-neighbours`, default=10), whether to add a `Odd ratio weighting` (default: off) and the parameter to be used for the node `Size` (default: cumFC).

edge threshold:

| 0 | 0.3 | | 1 |

0  0.1  0.2  0.3  0.4  0.5  0.6  0.7  0.8  0.9  1

N-neighbours:

○ 5  ● 10  ○ 25  ○ 100

☐ Odd ratio weighting

Size:

○ FC  ● cumFC  ○ centrality

**Cumulative enrichment**

In this plot, gene sets that are frequently shared with high enrichment will show a higher cumulative scores. You can choose between signed or `Absolute foldchange` in the options and use `order by` to order the plot by FC or cumulative FC (cumFC, default).

These plots are complementary to the plots provided by the **FC correlation** tab, to explore more in detail the nature of the correlation between contrasts.

CHAPTER

# NINETEEN

# SYSTEMSBIO

The final module of the platform is divided into three submodules: **Drug connectivity**, **Cell profiling** and **WGCNA**.

## 19.1 Drug connectivity

In the **Drug connectivity** submodule, users can correlate their signature with more than 5000 known drug profiles from the L1000 database, as well as with drug sensitivity profiles from the CTRP v2 and GDSC databases. Additionally, a separate list of shRNA- and cDNA-perturebed datasets from the L1000 database is also available (gene/L1000).

An activation-heatmap compares drug activation profiles across multiple contrasts. This facilitates to quickly see and detect the similarities between contrasts for certain drugs.

### 19.1.1 Settings panel

In the **Settings** panel, users can specify the contrast of their interest with the `Contrast` setting. Under `Analysis type` users can select from four databases, including the L1000 drug connectivity map (L1000/activity), the L1000 gene perturbation (L1000/gene) database, the CTRP v2 drug sensitivity (CTRP_v2/sensitivity) database and the GDSC drug sensitivity (GDSC/sensitivity) database (default: L1000/activity). The `only annotated drugs` option is used to exclude drugs without a known mechanism of action.

## 19.1.2 Drug enrichment

There are four main panels in the **Drug enrichment** tab:

**Drug connectivity**
> The Drug Connectivity panel correlates your signature with profiles from the L1000 (activity/L1000 and gene/L1000), CTRP and GDSC databases. It shows the top N=12 similar and opposite profiles as GSEA plots by running the GSEA algorithm on the contrast-drug profile correlation space.

**Enrichment table**
> Enrichment is calculated by correlating your signature with the profiles from the chosen database. Because of multiple perturbation experiments for a single small molecule, they are scored by running the GSEA algorithm on the contrast-small molecule profile correlation space. In this way, we obtain a single score for multiple profiles of a single small molecule. The table can be customised via the table *Settings* to only show annotated drugs.

**Mechanism of action**
> This plot visualizes the mechanism of action (MOA) across the enriched drug profiles. On the vertical axis, the number of drugs with the same MOA are plotted. You can switch to visualize between MOA or target gene. Under the plots *Settings*, users can select the plot type of MOA analysis: by class description (`drug class`) or by target gene (`target gene`). They can also apply q-value weighting for NES scoe values (`q-weighting`).

Plot type:
- ● drug class   ○ target gene
- ☐ q-weighting

**Activation matrix**
> The **Activation matrix** visualizes the correlation of small molecule profiles with all available pairwise comparisons. The size of the circles correspond to the strength of their correlation, and are colored according to their positive (red) or negative (blue) correlation to the contrast profile. The matrix can be normalised via the settings icon by ticking the `normalize activation matrix` option.

☐ normalize activation matrix

This tab can have many applications, which include understanding the MOA of a novel compund, identifying drugs that can be repurposed for treating a disease, identifying suitable partner drugs for the tested compound or target genes for intervention.

## 19.2 Cell Profiling

The **Cell Profiling** tab is specifically developed for the analysis and visualization of single-cell datasets. The main applications are identification of immune cell types and visualisations of markers, phenotypes, and proportions across the cells.

The **Cell type** tab infers the type of cells using computational deconvolution methods and reference datasets from the literature.

The **Mapping** tab provides a visualization of the inferred cell types matched to the phenotype variable of the data set, as well as a proportion plot visualizing the interrelationships between two categorical variables (so-called cross tabulation). This can be used to study the composition of a sample by cell type, for example.

The **Markers** tab provides potential marker genes, which are the top genes with the highest standard deviation within the expression data across the samples. It also generates a plot mimicking the scatter plots used for gating in flow cytometry analysis.

## 19.2.1 Settings panel

Users can filter relevant samples in the `Filter samples` settings under the the main `Options` in the input panel. They can also specify the `layout` for the figures by chooisng between pca, tsne or umap options (default: tsne).

## 19.2.2 Cell type

The **Cell type** tab contains two panels: **Cell type profiling** and **Phenotypes**.

**Cell type profiling** infers the type of cells using computational deconvolution methods and reference datasets from the literature. In the plot settings menu, users can select the reference dataset and the method for the cell type prediction in the `reference` and `method` settings, respectively. Currently, we have implemented a total of 7 methods (EPIC, DeconRNAseq, DCQ, I-NNLS, NNLM, correlation-based and a meta-method) and 9 reference datasets to predict immune cell types (4 datasets: LM22, ImmProt, DICE and ImmunoStates), tissue types (2 datasets: HPA and GTEx), cell lines (2 datasets: HPA and CCLE) and cancer types (1 dataset: CCLE). Not all methods or databases may be available for a dataset, the availability depends on the pre-processing done. From the settings, users can also sort plots by either probability or name and change the layout (`sort by`).

The **Phenotypes** tab displays plots that show the distribution of the phenotypes superposed on the t-SNE clustering. Often, we can expect the t-SNE distribution to be driven by the particular phenotype that is controlled by the experimental condition or unwanted batch effects. Users can customise the plot via the settings icon, where they can `label` the plot groups or add a legend instead.

The cell type profiling tab displays the two panels side by side.

## 19.2.3 Mapping

The **Mapping** panel contains two panels. The **Cell type mapping** panel contains a plot representing the cell type mapping across all samples. This plot can be customised via the *Settings* menu. Through it, users can change the `plot type` between a dotmap and a heatmap and select the `reference` dataset, select the analysis `method`. The reference datasets and the methods available are the same as indicated in the **Cell type profiling** panel under the **Cell type** tab. Users can also use `group by` to group samples by input phenotypes.

<div align="center">

plot type:

| dotmap ▾ |

reference:

| Immune cell (LM22) ▾ |

method:

| meta.prod ▾ |

group by:

| Sample_ID ▾ |

</div>

The **Proportions** panel contains a proportion plot visualizes the overlap between two categorical variables. This may be useful for bulk RNA datasets, as it can provide information about the proportion of different cell types in the samples. From the settings icon, users can select whwther to display the <cell type> (based on the chosen reference dataset) or select one of the available phenotypes on the x- and y-axes of the plot. By selecting a gene with `gene` they can also add an expression barplot that indicates the expression level (high or low is based on average sample expression for the gene) of the selected gene for each of the sample groups as well as adding the total number of read counts of the selected gene per sample group.

<div align="center">

x-axis:

| <cell type> ▾ |

y-axis:

| <cell type> ▾ |

gene:

| A1BG ▾ |

</div>

The two panels are displayed side by side in the tab.

## 19.2.4 Markers

The **Markers** tab consists of two panels: **Expression of marker genes** and **Cytometry plot**.

**Expression of marker genes** consists of 25 t-SNE plots of the genes with the highest standard deviation that could represent potential biomarkers. The red color shading is proportional to the (absolute) expression of the gene in corresponding samples. In the settings icon, users can specify the `Level` of the marker analysis: gene or gene set level. They can also restrict the analysis by selecting a particular functional group in the `Feature set`, where genes are divided into 89 groups, such as chemokines, transcription factors, genes involved in immune checkpoint inhibition, and so on (default: CD molecules (HGNC)). In addition, it is possible to filter markers by a specific keywords in the `Filter` setting and sort them by intensity (default) or name (`sort by`).

For each gene pairs combination, the panel also generates a cytometry-like plot (**Cyto plot**) of samples. The aim of this feature is to observe the distribution of samples in relation to the selected gene pairs. For instance, when applied to single-cell sequencing data from immunological cells, it can mimic flow cytometry analysis and distinguish T helper cells from other T cells by selecting the CD4 and CD8 gene combination. Under the plot settings icon, users can select

their prefered genes on the x- and y-axes in the `x-axis` and `y-axis`, respectively. They can also set the maximum number of bins for histgram distribution (`nbins`). This will be used to calculate the density distribution of the gene pairs selected in the `x-axis` and `y-axis`.



The two panels are displayed side by side in the tab.



## 19.3 WGCNA

The final submodule under **SystemsBio** is dedicated to weighted correlation network analysis (**WGCNA**), which serves the purpose of identifying clusters (modules) comprising highly correlated genes. These clusters can be summarized using either the module eigengene or an intramodular hub gene. WGCNA also facilitates the association of modules with each other and external sample traits through eigengene network methodology. Furthermore, it allows for the computation of module membership measures.

## 19.3.1 Settings panel

WGCNA modules are selected from the **Settings** panel (`select module`) and plots can be recalcluated based on the selected module. Under *Options*, the number of genes (`Number genes`, default=1000), the miinum module size (`Min. module size`, default=30), the `Power` (default=6), the `deepsplit` (default=2) and the `Merge cut height`, default=0.25) can be set.



## 19.3.2 WGCNA

The **WGCNA** tab consists of fiive panels (from left to right and top to bottom): **Gene dendrogram and gene modules**, **Scale independence and mean connectivity**, **TOM heatmap**, **Gene clustering** and **Module graph**.

**Gene dendrogram and gene modules**

In this panel, gene modules are detected as branches of the resulting cluster tree using the dynamic branch cutting approach. Genes inside a given module are summarized with the module eigengene. The module eigengene of a given module is defined as the first principal component of the standardized expression profiles.

**Scale independence and mean connectivity**

This panel is used for the the analysis of network topology for various soft-thresholding powers. The left plot shows the scale-free fit index (y-axis) as a function of the soft-thresholding power (x-axis). The right plot displays the mean connectivity (degree, y-axis) as a function of the soft-thresholding power (x-axis).

**TOM heatmap**

The panel displays the Topological Overlap Matrix (TOM) as a heatmap, which shows the correlation among gene module memberships.

**Gene clustering**

Dimensionality reduction maps colored by WGCNA module. Via the settings icon, the layout can be changed between tsne (default), pca and umap.



**Module graph**

The final panel contains the raph network of WGCNA modules, which represents the relationship betweem of the gene modules.



## 19.3.3 Modules

The **Modules** tab contains five panels (left to right, top to bottom): **Module-Trait relationships**, **Correlation network**, **Module Enrichment (plot)**, **Module genes** and **Module enrichment (table)**.

**Module-Trait relationships**

Module-trait analysis identifies modules that are significantly associated with the measured clinical traits by quantifying the association as the correlation of the eigengenes with external traits. The relationships between the various WGCNA modules and the phenotypic groups in the dataset are displayed as a heatmap, with shades of red indicating a negative correlation and shades of green

indicating a positive correlation. The continuous variables can be binarised via the settings icon (`binarize continuous vars`).

☐ binarize continuous vars

### Correlation network

A partial correlation graph centered on module eigengene with top most correlated features. Green edges correspond to positive (partial) correlation, red edges to negative (partial) correlation. Width of the edges is proportional to the correlation strength of the gene pair. The regularized partial correlation matrix is computed using the 'graphical lasso' (Glasso) with BIC model selection.

### Module Enrichment Plot

A plot that displays the functional enrichment of top most enriched genesets from a selection of GO and MSigDB genesets.

### Module genes

A table showing the genes in the WGCNA module selected via the **Settings** panel. The value me.rho represents the correlation between the gene expression and the module.

### Module Enrichment Table

Functional enrichment of the module calculated using Fisher's exact test. In this table, users can check mean expression values of features across the conditions for the selected module.

## 19.3.4 Eigengenes

The **Eigengenes** tab is used to visualise the network of eigengenes and study the relationships among the found modules. One can use the eigengenes as represetative profiles and quantify module similarity by eigengene correlation. For each module, we also define a quantitative measure of 'module membership' (MM) as the correlation of the module eigengene and the gene expression profile. This allows us to quantify the similarity of all genes to every module.

The tab contains two panels: **Eigengene clustering** and **Module membership (eigengene correlation)**.
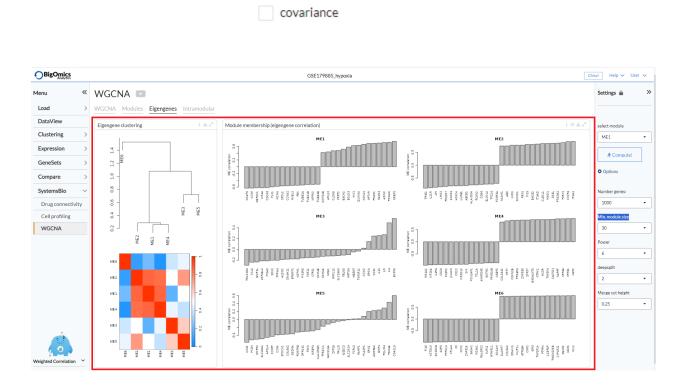
> **Eigengene clustering**
>> A cluster heatmap that shows the relationship between the different modules produced by the platform.

> **Module membership (eigengene correlation)**
>> This panels contains a series of plots for each one of the generated modules. For each module, we define a quantitative measure of module membership (MM) as the correlation of the module eigengene and the gene expression profile. This allows us to quantify the similarity of all genes on the array to every module. Users can also select to include the `covariance` for each gene (default:off).





## 19.3.5 Intramodular

The **Intramodular** tab is used to quantify associations of individual genes with the trait of interest (weight) by defining Gene Significance (GS) as (the absolute value of) the correlation between the gene and the trait. For each module, a quantitative measure of module membership (MM) as the correlation of the module eigengene and the gene expression profile is also defined. Using the GS and MM measures, users can identify genes that have a high significance for weight as well as high module membership in interesting modules.

The tab contains two panels: **Membership-trait heatmap** and **Membership vs. trait correlation**.

> **Membership-trait heatmap**

For each module, a quantitative measure of module membership (MM) as the correlation of the module eigengene and the gene expression profile is defined. This allows us to quantify the similarity of all genes on the array to every module and represent them as a heatmap.

### Membership vs. trait correlation

In this panel, the MM and trai correlations are represented as a series of scatterplots.

# REANALYZING PUBLIC DATASETS

To illustrate the use case of the Omics Playground, we reanalyzed different types of publics datasets, including microarray, bulk RNA-seq, single-cell RNA-seq and proteomic datasets to recapitulate the results.

## 20.1 Single-cell RNA-seq data

For single-cell RNA-seq data, we downloaded the melanoma data set GSE72056 of Tirosh et al.. Our platform recapitulates well the original findings of the paper. The t-SNE clustering (*Figure 1*) separates the different cell types. *Figure 2* and *Figure 3* show the volcano plot, MA plot and most differentially expressed genes between malignant and non-malignant cells. The CNV map (*Figure 4*) confirms the major chromosomal copy number variations found in the malignant cells. *Figure 5* shows high enrichment of a immune checkpoint signature, particularly concentrated in the T cells. The biomarker heatmap (*Figure 6*) highlights the marker genes for each cell type. Each gene cluster is furthermore automatically annotated with the most correlated gene sets (*Figure 7*).
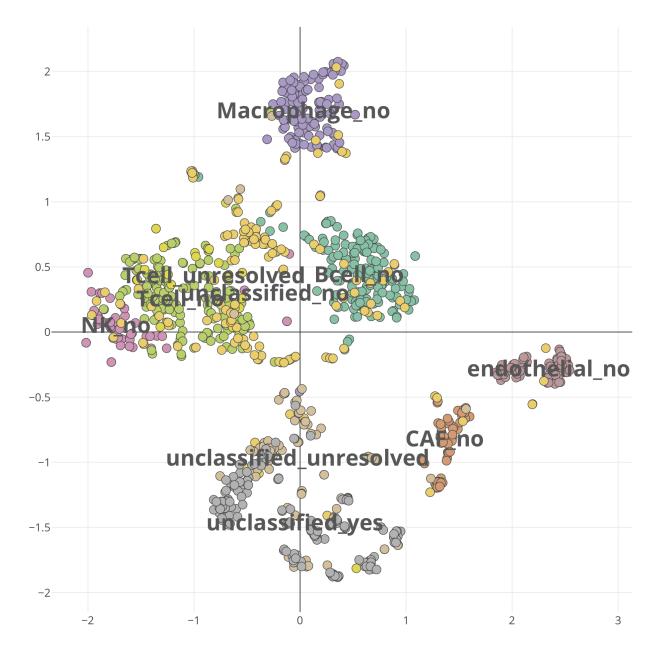
### 20.1.1 tSNE plot

**Figure 1**. The t-SNE clustering with cell type annotation for the `GSE72056-scmelanoma` dataset. To reproduce the figure on the platform, select and load `GSE72056-scmelanoma` dataset, and go to the **PCA/tSNE** panel of the **Clustering** module. From the plot *Settings*, set the `color:  group`, `layout:  tsne`, and leave other settings as default.
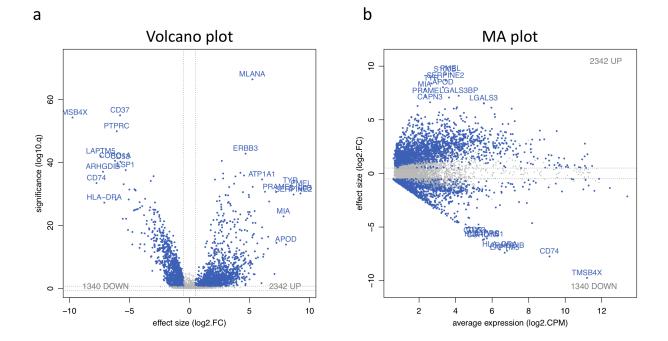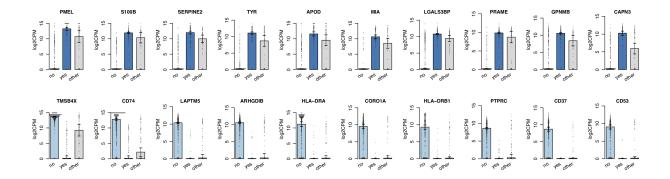
### 20.1.2 Volcano and MA plot

**Figure 2**. Volcano and MA plot for the malignant versus non-malignant contrast. To replicate the figure on the platform, go to the **Plots** panel of the **Expression** module. From the input slider, set the `Contrast:  yes_vs_no` and `Gene family:  all`.

### 20.1.3 Differentially expressed genes

**Figure 3**. Barplot of corresponding differentially expressed genes. To obtain the figure on the platform, go to the **Top genes** panel of the **Expression** module. From the input slider, set the `Contrast:  yes_vs_no` and `Gene family: all`.

a
## Volcano plot



b
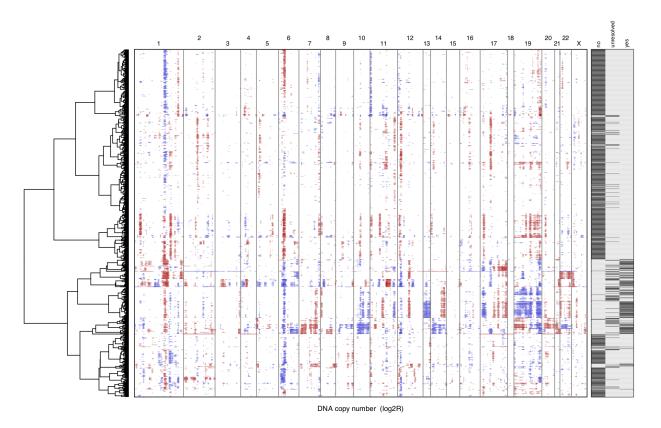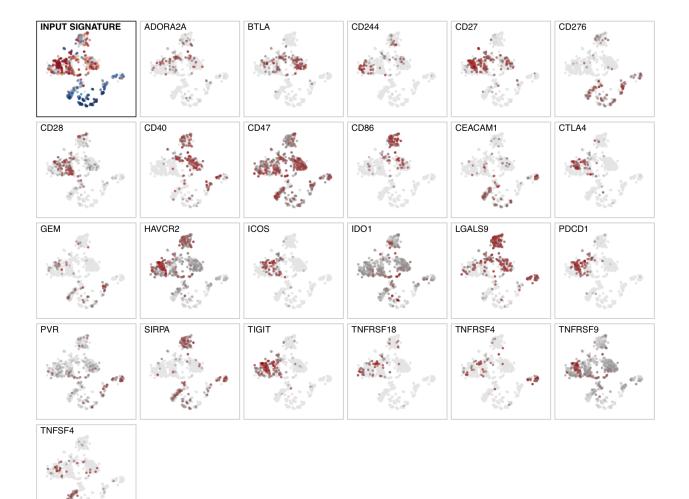## MA plot

## 20.1.4 Inferred copy number



**Figure 4**. Inferred copy number for sample Cy80. To reproduce the figure on the platform, go to the **CNV** panel of the **scProfiling** module. From the plot *Settings*, set the `Annotate with:  malignant` and `Order samples by: clust`.
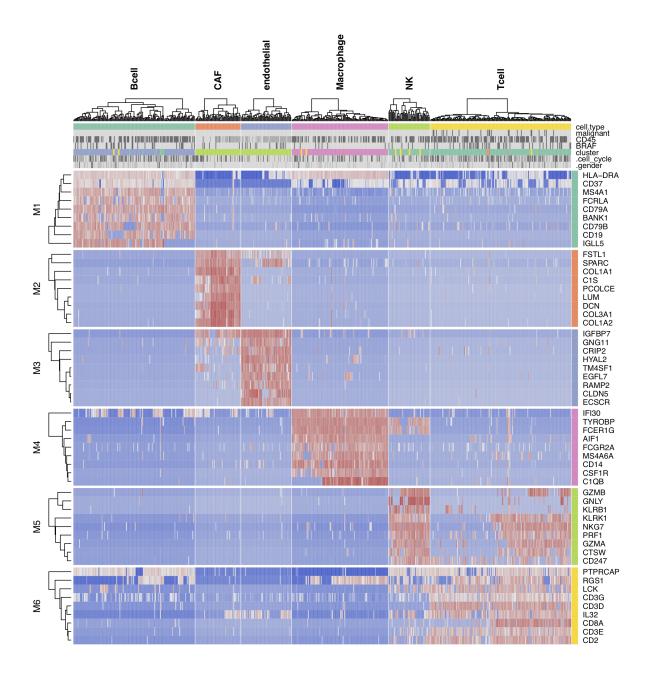
## 20.1.5 Immune checkpoint signature

**Figure 5**. Enrichment distribution for an immune checkpoint signature showing high enrichment in T and B cells . To regenerate the figure on the platform, go to the **Marker** panel in the **Signature** module. From the input slider, select `Contrast:  custom` and `Signature:  immune_chkpt` as it is provided in the sample list.

## 20.1.6 Biomarker heatmap

**Figure 6**. Biomarker heatmap for non-malignant cells. To reproduce the figure on the platform, go to the **Heatmap** panel in the **Clustering** module. From the input slider, set the `Filter samples:  cell.type={Bcell, CAF, endothelial, Macrophage, NK, Tcell}`. In the plot *Settings*, set `Plot type:  ComplexHeatmap`, `split by: cell.type`, and `top mode:  specific`.
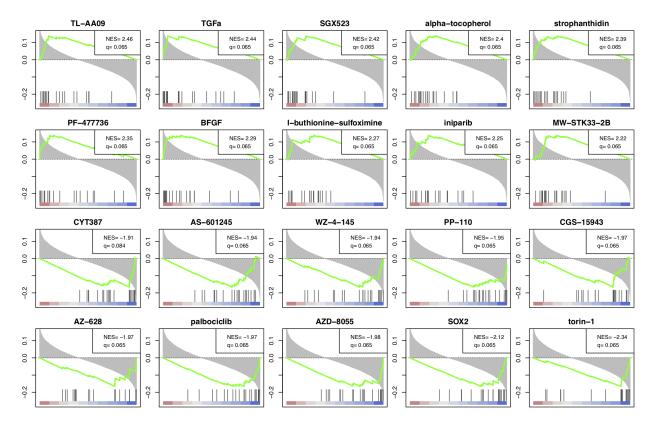
## 20.1.7 Annotate heatmap clusters



**Figure 7**. Enrichment annotation of corresponding heatmap clusters from the *Figure 6*. To reproduce the figure on the platform, generate the heatmap in *Figure 6* first, then go to the **Annotate clusters** panel. From the plot *Settings*, set the `Reference set: GOBP`.

## 20.2 Bulk RNA-seq Data

To elucidate the mechanism of action of a new drug, or for the intention of drug repurposing, it is often useful to find other drugs that have similar or opposing signatures compared to some given fold change profile. This example is illustrated using the RNA-sequencing dataset GSE114716 from Goswami et al., which contains CD4 T cells following ipilimumab therapy. *Figure 8* shows the top ranked drugs with most similar or most opposing signatures to ipilimumab, a novel monoclonal antibody targeting CTLA-4 used in tumour therapy to stimulate the immune system. The complete list contains several compounds that stimulate the immune system, such as alpha-tocopherol (Morel et al.), but also highlights compounds that are not commonly associated with the modulation of immune responses, such as strophanthidin, an intropic drug that has recently been shown to display pro-inflammatory activities (Karas et al.).

## 20.2.1 Drug enrichment profiles



**Figure 8**. Drug enrichment profiles for most similar and opposing drugs compared to ipilimumab treatment. To reobtain the figure on the platform, select and load `GSE114716-ipilimumab` dataset, go to the **Drug CMap** panel under the **Functional** module, and set the `Contrast:  Ipi_vs_baseline` from the plot *Settings*.

# 20.3 Microarray Data

In this section, we perform the heatmap clustering, biomarker selection and survival analysis using the GSE10846 from Lenz et al., which is the microarray gene expression dataset of diffuse large B-cell lymphoma (DLBCL) patients. *Figure 9* shows a hierarchical cluster heatmap of microarray gene expression data. *Figure 10* and *Figure 11* show the variable importance plot and a survival tree on the overall survival of the DLBCL patients, respectively.

## 20.3.1 Hierarchical cluster heatmap

**Figure 9**. Hierarchical cluster heatmap for `GSE10846-dlbcl` dataset. To replicate the figure, select and load `GSE10846-dlbcl` dataset on the platform. Go to the **Heatmap** panel of the **Clustering** module, and set the `Level: gene` and `Features:  all` from the input panel. In the plot *Settings*, set the `split by:  none` and `Top mode: pca`.

## 20.3.2 Variable importance plot



**Figure 10**. Variable importance plot. To replicate the figure, go to the **Biomarker** module, and set the `Predicted target: dlbcl.type` from the input panel.

## 20.3.3 Survival tree

**Figure 11**. Survival tree analysis for `GSE10846-dlbcl` dataset. To redproduce similar figures, go to the **Biomarker** module, and set the `Predicted target: OS.survival` from the input panel. Note that the survival tree is stochastically built up with some of the top features shown in *Figure 10*; Therefore, users can get a slightly different survival tree every time.

# 20.4 Proteomic Data

With larger data sets, often the number of contrasts increases and complicates the overall analysis. For example, the proteomics data set of Rieckmann et al. 2017 comprises 26 populations of seven major immune cell types, measured during resting and activated states. There are more than 300 possible comparisons to make. To gain a better overview, gene set activation matrix (*Figure 12*) help visualize the similarities between multiple contrasts on a functional level. Alternatively, similarities can be visualized as a connectivity graph (*Figure 13*). For the same data set, *Figure 14* shows a computed partition tree that classifies the major cell types.

Another example dataset is from Geiger et al., where the proteome profiles of activated vs resting human naive T cells at different times were compared. *Figure 15* shows the volcano plots corresponding to eight different statistical tests comparing time-dependent activation of T cells at 48h vs. 12h. We see that both standard t-test and the Welch t-test show much less power to detect significant genes compared to the other methods. The result from edgeR-QLF is close to those of the two limma based methods, while edgeR-LRT is very similar to the results of DESeq2-Wald.

### 20.4.1 Activation matrix

**Figure 12**. Gene Ontology activation matrix. To replicate the figure, select and load the `rieckmann2017-immprot` dataset, and go to the **GO graph** panel of the **Functional** module with default settings.

### 20.4.2 Contrast heatmap

**Figure 13**. Contrast heatmap for the `rieckmann2017-immprot` dataset. To generate the figure on the platform, go to the **Contrast heatmap** panel of the **Intersection** module with default settings.

### 20.4.3 Classification tree

**Figure 14**. Classification tree for the `rieckmann2017-immprot` dataset. To reproduce similar figures, go to the **Biomarker** module, and set the `Predicted target:  cell.type` from the input panel. Note that the classification tree is stochastically built up with some of the top features shown in *Figure 10*; Therefore, users can get a slightly different survival tree every time.

### 20.4.4 Volcano plots of methods

**Figure 15**. Volcano plots corresponding to eight different statistical methods comparing time-dependent expression of T cell activation at 48h vs. 12h. To regenerate the figure, select and load `geiger2016-arginine` dataset. Go to the **Volcano (methods)** panel under the **Expression** module, and set the `Contrast:  act48h_vs_act12h`.

# METHODS

Below are snippets that you can use to describe the methods when using the Omics Playground. These are just examples and you need to extract and modify the parts you used and need.

## 21.1 Batch correction

Batch effects, or contamination by unwanted variables, was identified by an F-test for the first three principal components. Continuous variables were dichotomized into high/low before testing. Highly confounding variables would appear as having high relative contribution in the first or second principal component, often higher than the variable of interest. Batch effects were also visually assessed (before and after correction) using annotated heatmaps and t-SNE plots colored by variables.

Batch correction was performed for explicit batch variables or unwanted covariates. Parameters with a correlation r>0.3 with any of variables of interest (i.e. the model parameters) were omitted from the regression. Correction was performed by regressing out the covariate using the 'removeBatchEffect' function in the limma R/Bioconductor package.

Technical correction was performed for intrinsic technical parameters such as library size (i.e. total counts), mitochondrial and ribosomal proportions, cell cycle and gender. These parameters were estimated from the data. The cell cycle was estimated using the Seurat R/Bioconductor package. Gender (if not given) was estimated by checking the absence/presence of expression of gender specific genes on the X/Y chromosome. Parameters with a correlation r>0.3 with any of the model parameters were omitted from the regression. Correction was performed by regressing out the covariate using the 'removeBatchEffect' function in the limma R/Bioconductor package.

Unsupervised batch correction was performed using 'surrogate variable analysis' (SVA) (Leek 2007) by estimating the latent surrogate variables and regressing out using the 'removeBatchEffect' function in the limma R/Bioconductor package.

## 21.2 Clustering

Heatmaps were generated using the ComplexHeatmap R/Bioconductor package (Gu 2016) on scaled log-expression values (z-score) using euclidean distance and Ward linkage. The standard deviation was used to rank the genes for the reduced heatmaps.

T-distributed stochastic neighbor embedding (t-SNE) was computed using the top 1000 most varying genes, then reduced to 50 PCA dimensions before computing the t-SNE embedding. The perplexity heuristically set to 25% of the sample size or 30 at maximum, and 2 at minimum. Calculation was performed using the *Rtsne* R package.

Uniform Manifold Approximation and Projection (UMAP) was computed using the top 1000 most varying genes, then reduced to 50 PCA dimensions before computing the UMAP embedding. The number of neighbours was heuristically

set to 25% of the sample size or 30 at maximum, and 2 at minimum. Calculation was performed using the *uwot* R package.

Principal component analysis (PCA) was computed using the *irlba* R package.

## 21.3 Statistical testing

Multi-method statistical testing. For gene-level testing, statistical significance was assessed using three independent statistical methods: DESeq2 (Wald test), edgeR (QLF test) and limma-trend (Love 2014; Robinson 2010; Ritchie 2015). The maximum q-value of the three methods was taken as aggregate q-value, which corresponds to taking the intersection of significant genes from all three tests.

Statistical testing of differential enrichment of genesets was performed using an aggregation of multiple statistical methods: Fisher's exact test, fGSEA (Korotkevich 2019), Camera (Wu 2012) and GSVA/limma (Hanzelmann 2013, Ritchie 2015). The maximum q-value of the selected methods was taken as aggregate meta.q value, which corresponds to taking the intersection of significant genes from all tests. As each method uses different estimation parameters (NES for GSEA, odd-ratio for fisher, etc.) for the effect size, for consistency, we took the average log fold-change of the genes in the geneset as sentinel value. We used more than 50000 genesets from various public databases including: MSigDB (Subramanian 2005; Liberzon 2015), Gene Ontology (Ashburner 2000), and Kyoto Encyclopedia of Genes and Genomes (KEGG) (Kanehisa 2000).

## 21.4 Functional analysis

Graph-weighted GO analysis. The enrichment score of a GO term was defined as the sum of q-weighted average fold-changes, $(1-q)*logFC$, of the GO term and all its higher order terms along the shortest path to the root in the GO graph. The fold-change of a gene set was defined as the average of the fold-change values of its member genes. This graph-weighted enrichment score thus reflects the enrichment of a GO term with evidence that is corroborated by its parents in the GO graph and therefore provides a more robust estimate of enrichment. The activation map visualizes the scores of the top-ranked GO terms for multiple contrasts as a heatmap.

KEGG pathway visualization was performed using the Pathview R/Bioconductor package using the foldchange as node color.

## 21.5 Cell type profiling

Cell type profiling was performed using the LM22 signature matrix as reference data set (Chen 2018). We have evaluated a total of 6 computational deconvolution methods: DeconRNAseq (Gong 2013), DCQ (Altboum 2014), I-NNLS (Abbas 2009), NNLM (Lin 2020), rank-correlation and a meta-method. For NNLM, we repeated NNLM for non-logarithmic (NNLM.lin) and ranked signals (NNLM.rnk). The latter meta-methods, meta and meta.prod, summarize the predictions of all the other methods as the mean and/or geometric mean of the normalized prediction probabilities, respectively.

[1] Gong T, Szustakowski JD. DeconRNASeq: a statistical framework for deconvolution of heterogeneous tissue samples based on mRNA-Seq data. Bioinformatics. 2013.

[2] Altboum Z, et al. Digital cell quantification identifies global immune cell dynamics during influenza infection. Mol Syst Biol. 2014 Feb 28;10(2):720.

[3] Abbas A, et al. Deconvolution of Blood Microarray Data Identifies Cellular Activation Patterns in Systemic Lupus Erythematosus, PLOS One, 2009.

[4] Lin X, Boutros PC. Optimization and expansion of non-negative matrix factorization. BMC Bioinformatics. 2020.

[5] Chen B, et al. Profiling Tumor Infiltrating Immune Cells with CIBERSORT. Methods Mol Biol. 2018.

## 21.6 Scripting and visualization

Data preprocessing was performed using bespoke scripts using R (R Core Team 2013) and packages from Bioconductor (Huber 2015). Statistical computation and visualization have been performed using the Omics Playground version vX.X.X (Akhmedov 2020).

## 21.7 REFERENCES

Akhmedov M, Martinelli A, Geiger R and Kwee I. Omics Playground: A comprehensive self-service platform forvisualization, analytics and exploration of Big Omics Data. NAR Genomics and Bioinformatics, Volume 2, Issue 1, March 2020,

Ashburner et al. Gene ontology: tool for the unification of biology. Nat Genet. May 2000;25(1):25-9.

Huber W, et al. (2015) Orchestrating high-throughput genomic analysis with Bioconductor. Nature Methods 12:115-121; doi:10.1038/nmeth.3252

Kanehisa, M. and Goto, S.; KEGG: Kyoto Encyclopedia of Genes and Genomes. Nucleic Acids Res. 28, 27-30 (2000).

Leek J., Storey J. Capturing heterogeneity in gene expression studies by 'surrogate variable analysis' PLoS Genet. 2007

Love MI, Huber W, Anders S (2014). "Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2." Genome Biology, 15, 550.

R Core Team (2013). R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. URL http://www.R-project.org/.

Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK (2015). "limma powers differential expression analyses for RNA-sequencing and microarray studies." Nucleic Acids Research, 43(7)

Robinson MD, McCarthy DJ, Smyth GK (2010). "edgeR: a Bioconductor package for differential expression analysis of digital gene expression data." Bioinformatics, 26(1), 139-140.

# FREQUENTLY ASKED QUESTIONS (FAQ)

## 22.1 Can I use LIMMA, EdgeR and DESeq2 for my proteomics data?

EdgeR and DESeq2 are statistical methods based on the negative binomial model (an overdispersed Poisson model). Poisson and negative binomial models, naturally account for heteroscedasticity and zero values in the data. Also LIMMA is widely used in RNA-seq analysis, and is an emperical Bayesian method based on moderated t-test statistics. While these methods were originally conceived for the for differential expression analysis in RNA-seq data, there is increasingly more acceptance that these methods can also be used for proteomics data.

Langley and Mayer (2015) assessed seven methods for differential expression analysis in proteomics, and showed that DESeq to outperform the more commonly used t-test. Kammers et al. (2015) show that LIMMA shows better results than the t-test. Branson and Fretais (2016) also highlighted the statistical analogy of proteomics label-free spectral count quantification with count data from RNA sequencing and propose to use edgeR, DESEq and baySeq for proteomics data. Gregori et al. (2013) uses EdgeR and Poisson based methods for their LC/MS-MS data. They also created the 'msmsTests' R package. Medo et al (2019) use EdgeR for their study on missing values in the differential analysis of proteomic and phosphoproteomics data. Gatto proposes both edgeR and LIMMA based methods for proteomics statistical analysis. Chen et al. (2020) reviewed bioinformatics methods for proteomics data analysis and reaffirmed that LIMMA can achieve more robust and accurate results than the traditional t-test.

Notice that proper scaling/normalization of quantitative proteomics data is important before using RNA-seq based methods as sometimes the proteomics intensity values may be far greater than those found in RNA-seq. Therefore Omics Playground scales proteomics data automatically to 'counts per million' (CPM). We have also seen that batch correction (e.g. using ComBat) may improve the downstream statistical analysis in proteomics.

References:

1. Kammers et al. "Detecting significant changes in protein abundance". EuPA Open Proteomics Volume 7, June 2015.

2. Langley SR, Mayr M. "Comparative analysis of statistical methods used for detecting differential expression in label-free mass spectrometry proteomics". J Proteomics. 2015.

3. Branson OE, Freitas MA. "A multi-model statistical approach for proteomic spectral count quantitation". J Proteomics. 2016.

4. Gregori et al. "msmsTests-package: LC-MS/MS Differential Expression Tests". R package.

5. Gregori et al. "An Effect Size Filter Improves the Reproducibility in Spectral Counting-based Comparative Proteomics." Journal of Proteomics, 2013.

6. Medo, M, Aebersold, DM and Medová, M "ProtRank: bypassing the imputation of missing values in differential expression analysis of proteomic data." BMC Bioinformatics 20, 563 (2019).

7. Gatto L. "Bioconductor tools for mass spectrometry and proteomics", https://lgatto.github.io/bioc-ms-prot/lab.html#8_statistical_analysis

8. Chen et al. "Bioinformatics Methods for Mass Spectrometry-Based Proteomics Data Analysis." Int J Mol Sci. 2020;21(8):2873, 2020.

## 22.2 How are duplicated gene/protein names handled in the counts file?

Duplicated row identifiers (genes/proteins with same name) are handled by summing up their linear intensities/counts. If the data was in logarithm, it will be (hopefully) automatically detected and exponentiated. The rational of summing up the counts (or linear intensities in proteomics) is that we don't differentiate between possible gene/protein isoforms and sum them up as a group. If you want to retain the isoforms, you may keep the names as GENE.1 and GENE.2 but you must turn off any gene filter. However as currently such gene/protein variants are not recognized in the gene sets, this will result in wrong enrichment test.

## 22.3 How are "missing values" handled in the counts file?

At the moment missing values are imputed to 0 (zero). The Omics Playground uses the function *pgx.createPGX()* in the file pgx-compute.R. Note that only real NA (or empty) values in the counts file are supposed to be "missing". Zero values are assumed to be real zeros. If zero-value imputaton is not what you want, and you want to impute your "missing" values differently, (currently) you must do that manually before uploading the CSV file.

## 22.4 How do I switch between stable, beta & testing docker versions?

You can test the latest features of the Omics Playground by switching to a more experimental docker version. Omics Playground has three different docker versions:

- Stable channel (latest): This channel is fully tested by the Omics Playground team, and is the best choice to avoid crashes and other problems. It's updated roughly every 2–3 months for minor changes, and every 6 months for major changes.

- Beta channel (beta): To view upcoming improvements and features with low risk, use the Beta channel. It's updated roughly every month, with updates before the Stable channel gets them.

- Testing channel (testing): The Testing channel gets updated once or twice weekly and is intended only for developers. This build might have serious bugs or highly experimental features. The Testing channel always uses the latest source code and is build irregularly whenever needed.

You can switch between versions by pulling the docker image corresponding to the version latest, beta or testing.

# CONTACT

If you have further questions about how to use Omics Playground, please ask your question in our user forum at Google Groups.

If you've discovered a bug, have a feature request, or want to contribute to the code, please go to our GitHub repository.

If you have other questions you can email us.

# CITATION

To cite Omics Playground in publications please use:

**Akhmedov M, Martinelli A, Geiger R and Kwee I.** Omics Playground: A comprehensive self-service platform forvisualization, analytics and exploration of Big Omics Data. *NAR Genomics and Bioinformatics, Volume 2, Issue 1, March 2020, lqz019; doi:10.1093/nargab/lqz019*

# LICENSE

Omics Playground is open source: you are allowed to download the source code and make your own edits. Although Omics Playground is open source, the software is distributed under a dual license: for either non-commerical or commercial use.

## 25.1 Non-commercial

If you are a non-profit company, academic, or use our products for personal use, you may enjoy our software for free under a Creative Commons (CC) Attribution-NonCommercial licence.

## 25.2 Commercial/governmental

If you are not entitled to the non-commercial license, or you are a governmental agency, you must purchase a commercial license. Buying such a license is mandatory as soon as you develop commercial activities involving Omics Playground software. These activities include: offering paid services to customers as an ASP, SaaS, professional consulting, or shipping Omics Playground with a closed source product.