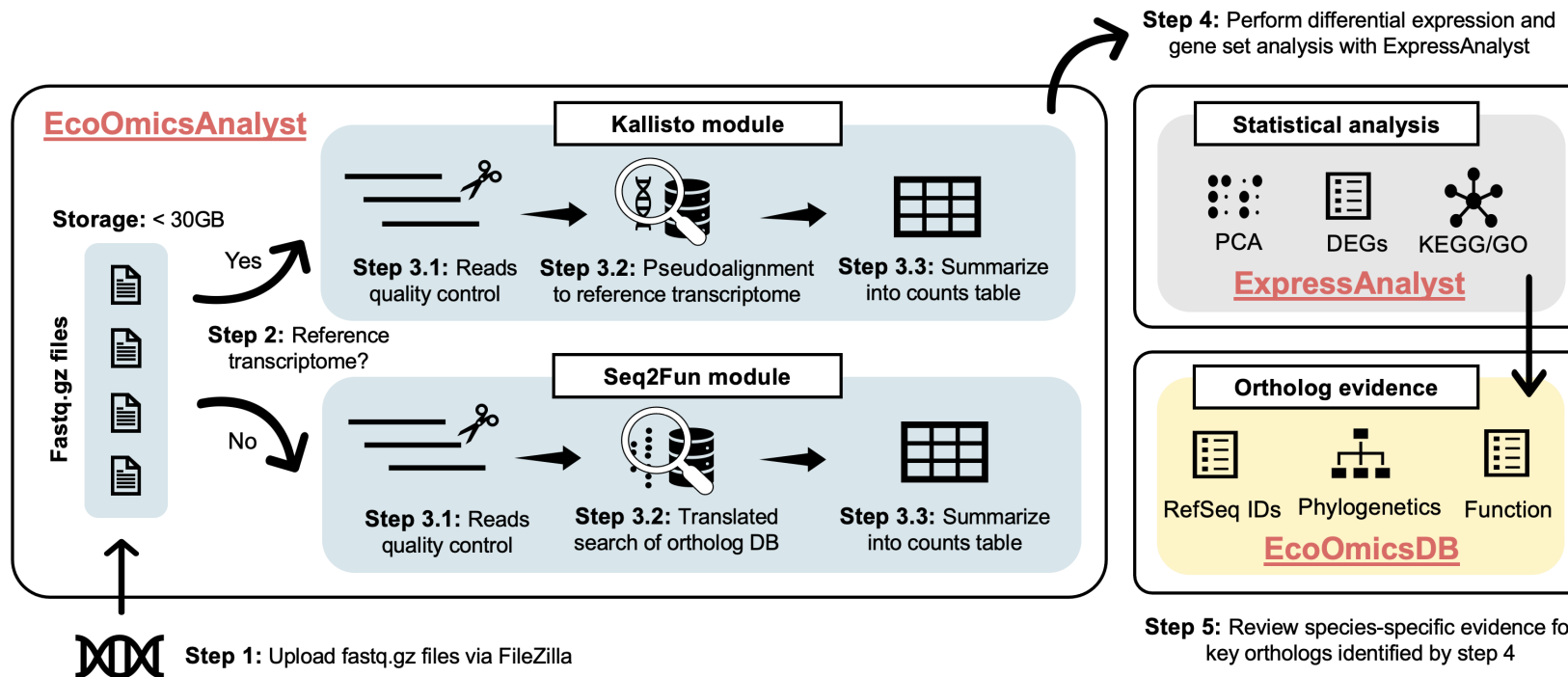


ExpressAnalyst for non-model organisms

-- analyzing Seq2Fun counts tables

Introduction

- ExpressAnalyst is part of our comprehensive solution for RNA-seq analysis for species without a reference genome:



Scope of this tutorial

- FASTQ files from three salamander species (2 had no reference genome) were previously processed with Seq2Fun using EcoOmicsAnalyst
 - See the 'Tutorials' tab at www.ecoomicsanalyst.ca
- Here, we demonstrate statistical and functional analysis of the resulting Seq2Fun ortholog counts table with ExpressAnalyst
- See tutorials I-IV for other ExpressAnalyst functions

Case study

Comparative transcriptomics of limb regeneration: Identification of conserved expression changes among three species of *Ambystoma*

Varun B. Dwaraka^{a,b,*}, Jeramiah J. Smith^a, M. Ryan Woodcock^c, S. Randal Voss^{b,d,e}

- RNA-seq measured in tissue samples collected after limb amputation
 - Immediately (time0) and after 24 hours (time24)
 - Three species (AND, MAC, MEX)
 - Three replicates per group: $(3 \times \text{time0} + 3 \times \text{time24}) \times 3 \text{ species} = 18 \text{ samples}$
 - One outlier removed (MEX, time0)
- Original study quantified MEX with reference transcriptome, AND and MAC with *de novo* transcriptome
 - Each species analyzed independently; results compared afterwards
- Here, data from all species are quantified with Seq2Fun
 - All data analyzed together; species accounted for during statistical analysis

Analysis objectives

- Determine similarities and differences in the transcriptomics of limb regeneration across three salamander species:
 - At whole transcriptome level (PCA plot)
 - In lists of differentially expressed genes
 - At functional level through pathway analysis

Computer and browser requirements

- A modern web browser with JavaScript enabled
- Supported browsers include Chrome, Safari, Firefox, and Internet Explorer 9+
- For best performance and visualization, use:
 - Latest version of Google Chrome
 - A computer with at least 4GB of physical RAM
 - A 15-inch screen or bigger (larger is better)
- Browser must be WebGL enabled for 3D scatter visualization
- 50MB limit for data upload
 - ~300 samples for gene expression data with 20 000 genes

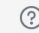
Expression Profiling



ExpressAnalyst


 Home

 Tutorials

 User Forum

 Updates

 Contact

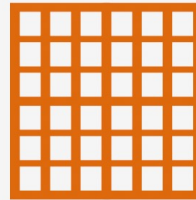
 About

Choose a module below to start analysis



- Starting from a list of gene IDs

Enrichment Analysis



- Starting from a single gene expression table

Expression Profiling



- Starting with several gene expression tables

Meta Analysis


Start here

Please use [OmicsForum](#) to ask questions related to ExpressAnalyst

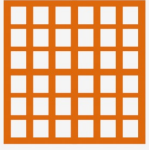
Find the data

Try our example data

<input type="radio"/> Estrogen	Affymetrix Human Genome U95 GeneChip (hgu95av2) data, normalized, log 2 scale (8 samples)	Gene expression of a breast-cancer cell line (source) . Estrogen Receptor (ER) : present, absent; Time (hour) : 10, 48
<input type="radio"/> Endotoxin	Illumina BeadArrays - Refseq ID, normalized, log 2 scale (12 samples)	Gene expression in human PBMC using LPS as inducer (details) Treatment : Control, LPS, LPS_LPS; Donor : 21, 46, 86, 92
<input type="radio"/> C. japonica toxicity	RNAseq data Entrez Gene ID, raw counts (15 samples)	Gene expression response in C. japonica from an early life stage toxicity experiment Treatment : Control, Medium, High;
<input checked="" type="radio"/> Non-model organisms	RNAseq data Seq2Fun ID, raw counts (17 samples)	Comparative transcriptomics of limb regeneration (details) Time : Time0, Time24; Species : <i>A. mexicanum</i> (MEX), <i>A. andersoni</i> (AND), <i>A. maculatum</i> (MAC).

 Submit

Click the 4th link to download the dataset



- Starting from a single gene expression table

Expression Profiling

Data Format: Single Matrix

Sample names

Meta-data


Gene IDs


#NAME	SRR7499348	SRR7499349	SRR7499351	SRR7499352	SRR7499353	SRR7499354	SRR7499355	SRR7499356	SRR7499357	SRR7499358	SRR7499359	SRR7499360	SRR7499361	SRR7499362	SRR7499363	SRR7499364	SRR7499365
#CLASS:Time	time0	time0	time0	time0	time0	time0	time0	time0	time24	time24	time24	time24	time24	time24	time24	time24	time24
#CLASS:Species	MEX	MEX	AND	AND	AND	MAC	MAC	MAC	MEX	MEX	MEX	AND	AND	AND	MAC	MAC	MAC
s2f_0000000100	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
s2f_0000000103	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
s2f_0000000105	13	3	2	1	6	8	6	10	5	7	5	2	6	4	4	0	4
s2f_0000000106	5	8	13	6	5	1	0	3	1	8	3	0	8	5	0	0	1
s2f_0000000107	49	45	41	63	44	22	26	41	49	52	51	40	46	47	3	9	42
s2f_0000000109	3	3	0	0	3	0	0	0	1	0	1	2	0	0	0	0	0
s2f_000000011	0	0	2	0	1	0	0	0	0	1	0	0	0	1	0	0	1
s2f_0000000110	0	0	0	0	1	0	2	3	0	0	1	1	0	0	0	0	3
s2f_0000000112	2	2	0	0	1	0	0	0	0	1	0	0	0	0	0	0	1
s2f_0000000115	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
s2f_0000000119	2	2	2	2	5	3	1	1	4	3	1	2	3	2	2	1	7
s2f_000000012	18	9	16	26	10	0	1	0	19	12	15	15	16	10	0	1	0
s2f_0000000120	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	0	0
s2f_0000000123	3	2	0	2	3	1	0	0	1	1	1	1	2	0	0	0	0
s2f_0000000124	2	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
s2f_0000000125	2	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	1
s2f_0000000128	0	1	0	0	0	0	0	1	2	0	0	0	0	0	0	0	0
s2f_0000000126	0	0	0	0	0	1	0	0	0	0	0	1	1	2	0	0	0
s2f_000000013	2	0	1	1	0	0	0	0	4	0	0	0	1	0	0	0	0
s2f_0000000130	0	0	1	0	1	0	0	0	0	2	3	2	0	0	0	0	1
s2f_0000000131	87	66	89	111	79	37	55	77	100	107	106	97	112	90	16	15	103

The data file can be tab delimited
(.tab/.txt) or comma delimited (.csv)

Data Upload and Annotation

Access tutorial, user forum, and updates pages

 > Upload > Download

 Navigate to:

Upload a gene expression table

ExpressAnalyst currently supports gene expression profiling and functional analysis for 25 organisms based on user feedback. including 11 model species, 5 pathogens and 9 ecological species. In addition, ExpressAnalyst also supports generic annotation based on KEGG orthologs (KO), as well as custom annotation. If your organism is not within the list, leave the **organism unspecified**, and you can still perform basic expression profiling such as differential analysis, volcano plot, heatmap clustering, etc.

Upload your gene expression table

Specify organism

Generic/Species independant

▼

Data type

Bulk RNA-seq data (counts)

▼

ID type

Seq2Fun Ortholog ID

▼

Gene-level summarization

Sum

▼

Data File

+ Choose

counts.txt 1.2 MB

Submit

1

Specify organism: Select "Generic" for Seq2Fun-processed data for species without reference genome

2

The gene level summarization depends on the data type. Microarrays produce intensity data so duplicate probes should be averaged (mean or median). RNAseq produce counts data, so multiple gene transcripts should be added (sum).

3

Specify ID type

4

Upload data file (for formatting see p. 5)

5

Click "Submit" then "Proceed"

<< Previous

>> Proceed

10

Quality Check: View Processing Results

[Home](#) > [Upload](#) > [Quality Check](#) > [Download](#)▼ Navigate to:

Data Quality Check

The uploaded samples are summarized below, together with several graphical outputs commonly used for quality check.

Data type:	RNA count table
Total feature number:	16755
Matched gene number:	16755 (100%)
Sample number:	17
Number of experimental factors:	2
Total read counts:	1.03e+08
Average counts per sample:	6.05e+06
Maximum counts per sample:	9.46e+06
Minimum counts per sample:	2.92e+06
Group names:	Two factors found - Time: time0; time24 Species: AND; MAC; MEX

Check the processing results to ensure correct sample size, experimental factors, and adequate gene annotation

Box plot

Count sum

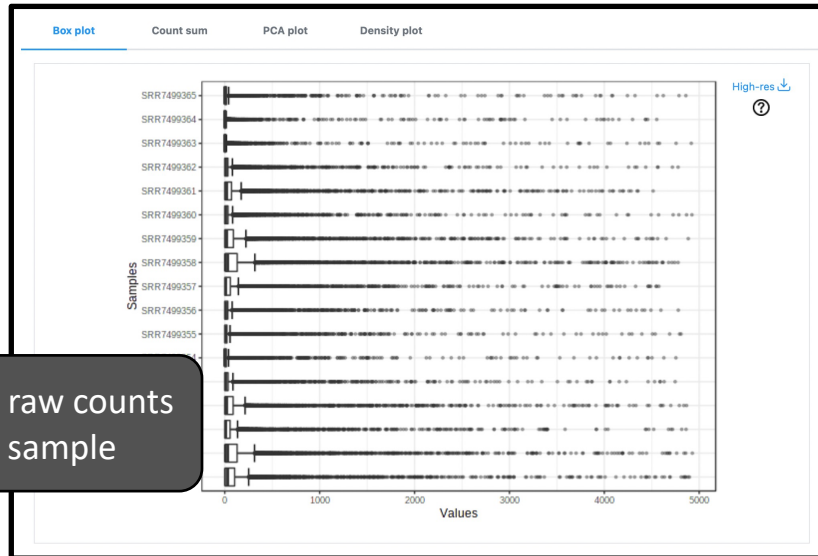
PCA plot

Density plot

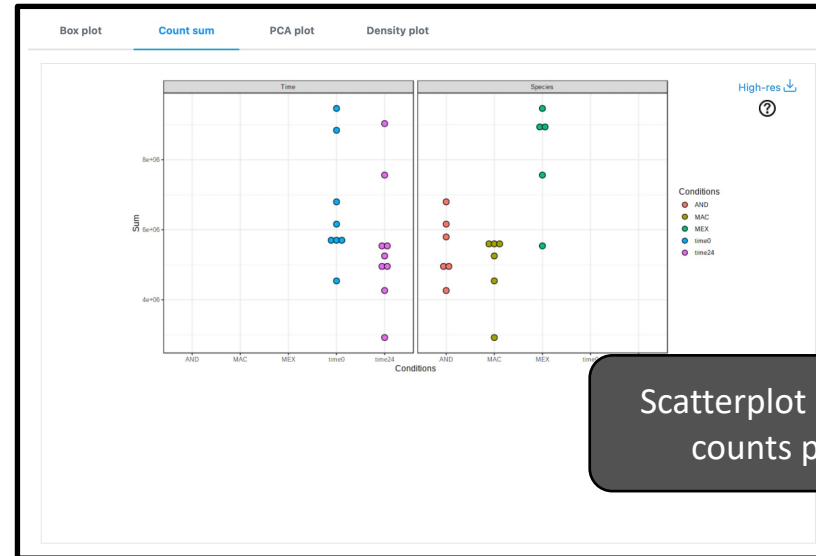
View common QA/QC plots to assess data quality

[<< Previous](#)[>> Proceed](#)

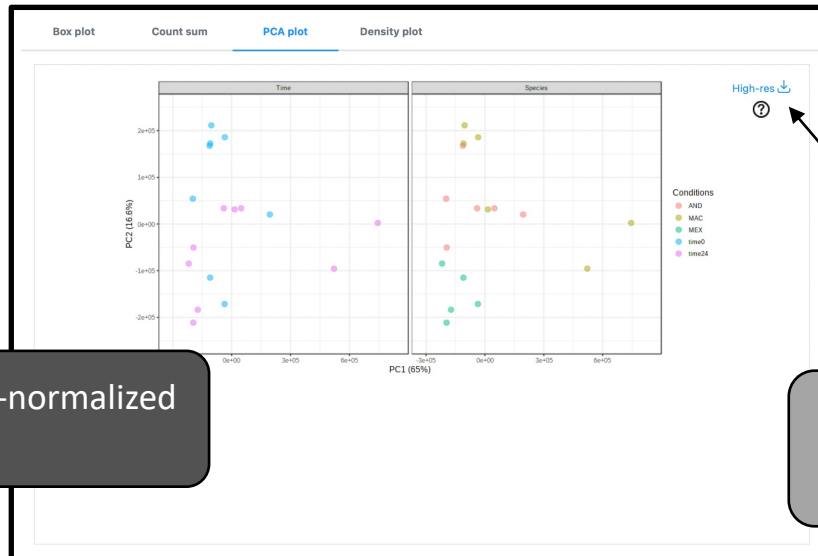
Quality Check: View Quality Control (QC) Plots



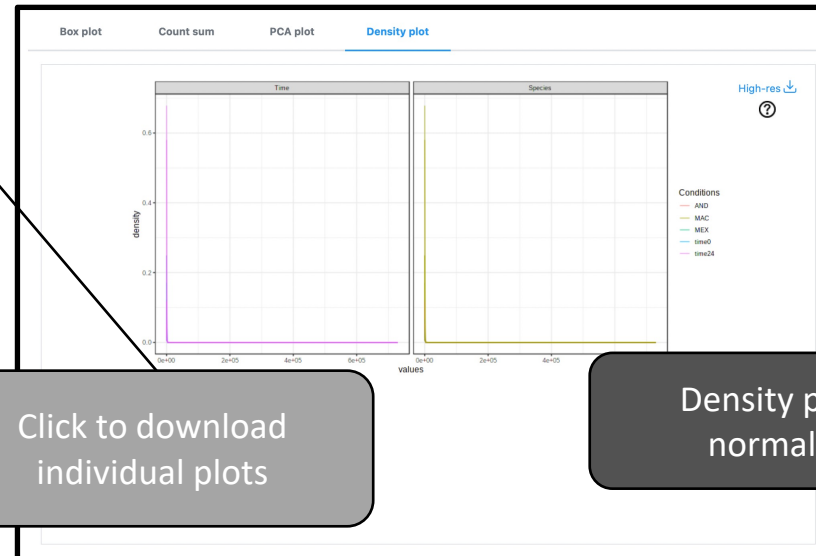
Boxplot showing raw counts per gene per sample



Scatterplot showing total counts per sample



PCA plot of non-normalized data



Density plot of non-normalized data

Click to download individual plots

Data Filtering and Normalization

Filtering increases statistical power by removing unresponsive genes prior to differential expression analysis (DEA). Proper normalization is essential to draw sound conclusions from the results of DEA.

[Home](#) > [Upload](#) > [Quality Check](#) > [Normalization](#) > [Download](#)

Data Filtering & Normalization

Filtering serves to remove data that are unlikely to be informative or simply erroneous. **Normalization** is crucial for a reliable detection of transcriptional differences, and to ensure that the expression distributions of each sample are similar across the entire experiment.

Filtering:

Variance filter:

15 ?

Low abundance:

4 ?

Filter unannotated genes:

☒

Normalization:

- ☐ None
☒ Log2-counts per million
☐ Upper Quantile Normalization
☐ Trimmed Mean of M-values
☐ Relative log expression normalisation

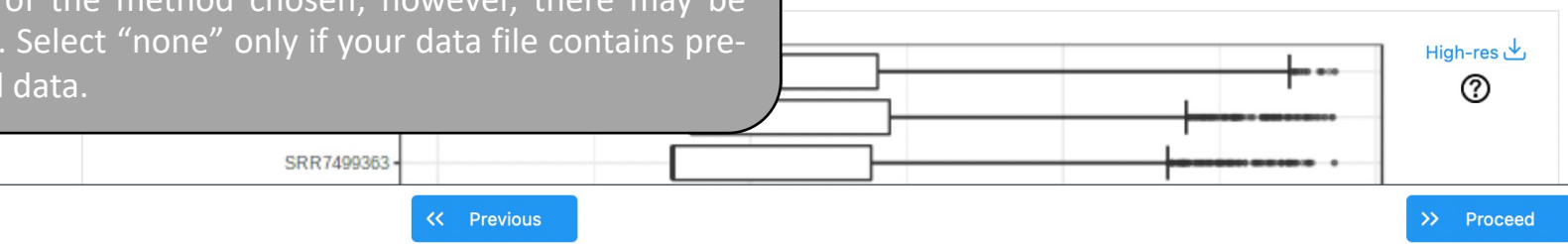
▶ Submit

1 Adjust variance and abundance filter: This determines how many genes are excluded from downstream analysis. The numbers represent percentiles – here, the 15th percentile of genes sorted by sample variance will be excluded; the 4th percentile of low count genes will be excluded.

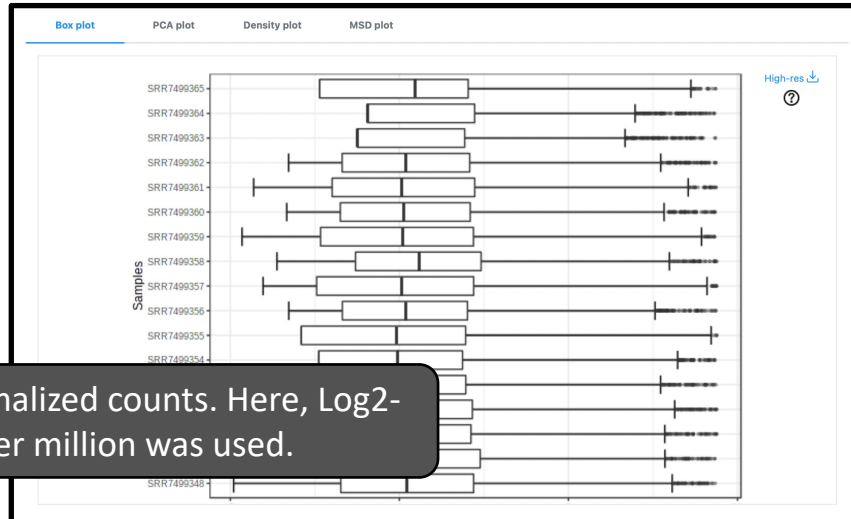
3 Click “Submit” to perform normalization. Effects will be reflected in the QC plots produced.

2 Select normalization method: All these are established methods commonly used for normalizing raw data. Downstream differential analysis results should be similar regardless of the method chosen; however, there may be exceptions. Select “none” only if your data file contains pre-normalized data.

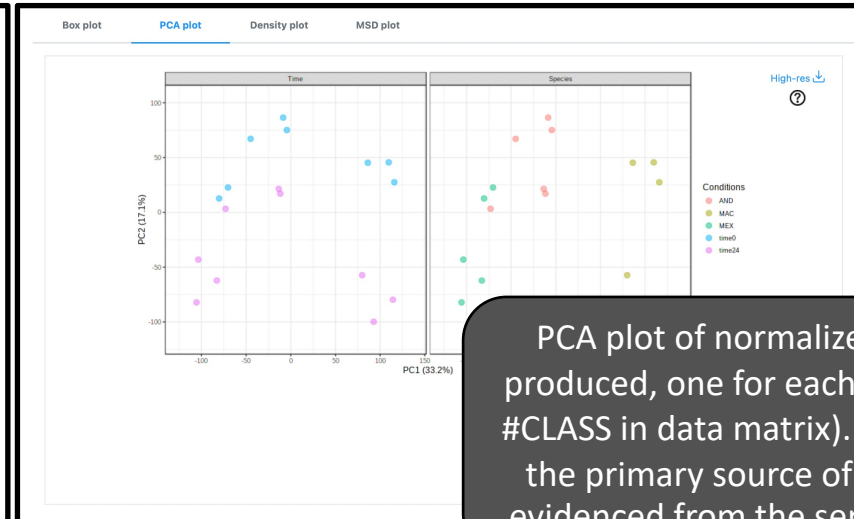
4 Click “Proceed” to continue analysis



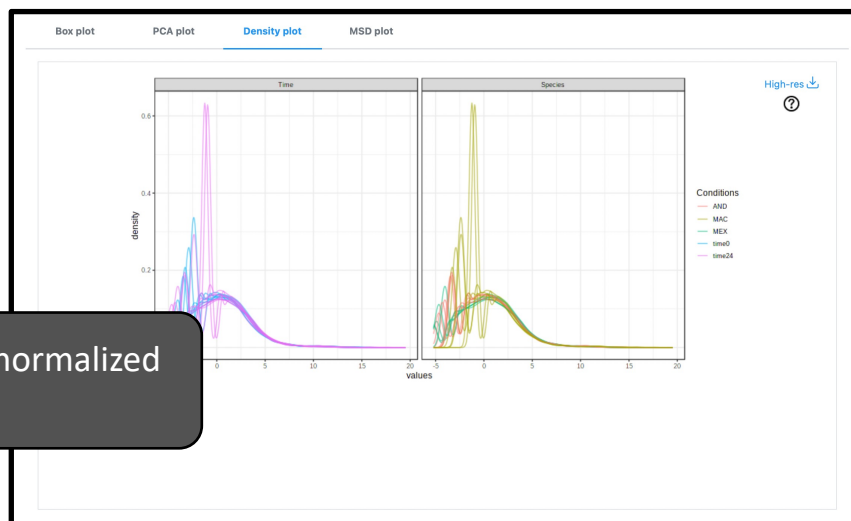
Data Filtering and Normalization



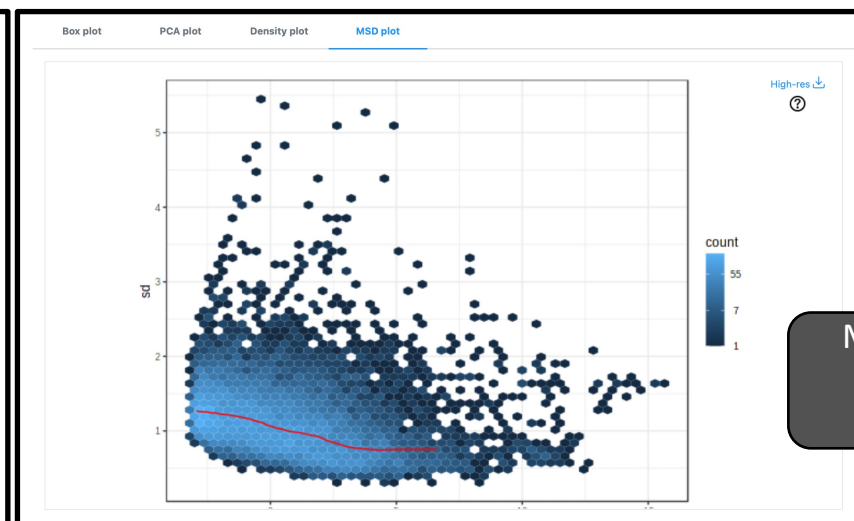
Box plot of normalized counts. Here, Log2-counts per million was used.



PCA plot of normalized data. Here, two plots are produced, one for each meta-data factor (defined by #CLASS in data matrix). Here we see that “Species” is the primary source of variability in the dataset, as evidenced from the separation of groups along PC1.



Density plot of normalized data



MSD (Mean vs. Standard Deviation) plot of normalized data

Perform DEA

The purpose of DEA is to group samples by a factor (ex: “time”, “treatment”, “sex”) and then perform statistical tests comparing groups within that factor for significant changes.

Home > Upload > Quality Check > Normalization > Differential Analysis > Sig. Genes > Download ▼ Navigate to:

Differential Expression Analysis

Statistical method

☒ Limma ☐ EdgeR ☐ DESeq2 ?

Study Design

Primary Factor: Time ?

Secondary Factor: Species ? This is a blocking factor ☒ ?

Comparison of Interest

☒ Specific comparison: time24 ? versus time0 ?

☐ Against a common control: time0 ?

☐ Nested comparisons: time0 vs. time24 ? versus time0 vs. time24 ? Interaction only ☒ ?

☐ Pairwise comparisons ?

⚙️ Submit

>> Proceed

2

Establish Study Design

Here, we had two factors in our meta-data: “Time”, and “Species.” We want to use “Time” as the Primary Factor to show the changes in gene expression over time following an event.

We selected “Species” as a “Blocking Factor” because this takes into account the variability attributed to species when calculating significance values for the Primary Factor (“Time”)

1

Select statistical method: All three methods are validated for DEA. Only “Limma” can be used for microarray data or for RNAseq data that has been normalized prior to being uploaded.

3

Establish Comparison of Interest (ex: “Treatment” vs. “Control”)

4

Click “Submit” to perform DEA.

5

Click “Proceed” to continue analysis

View Differentially Expressed Genes (DEGs)

Home > Differential Analysis > Upload > Quality Check > Normalization > Sig. Genes > Download

Select Significant Genes

Sig. Thresholds

Adjusted p-value: 0.1 ?

Log2 fold change: 0.0 ?

Submit

Total sig. genes: 2780

Download





Click "Submit" to filter by defined thresholds

Click "Download" retrieve gene expression list

Define significance thresholds: Modifying these values will yield different amount of DEGs

Here, we see 2780 DEGs based on our defined criteria for significance thresholds. Default values are Adjusted p-value: 0.05 and Log2 fold change: 1.

Gene ↑↓ View Details logFC ↑↓ AveExpr ↑↓ t ↑↓ P.Value ↑↓ adj.P.Val ↑↓ B ↑↓

PLEK2	 EODB	4.2201	3.221	14.125	3.3169E-12	4.7235E-8	17.66
MMP8	 EODB	6.2183	5.1073	12.795	2.1574E-11	1.0558E-7	15.975
TNC	 EODB	4.3384	4.4301	12.58	2.9597E-11	1.0558E-7	15.687
LAMB3	 EODB	4.6029	3.0128	12.579	2.9655E-11	1.0558E-7	15.685
MMP12	 EODB				8.1433E-11	2.2499E-7	14.757
MMP9	 EODB				1.0474E-10	2.2499E-7	14.524
KRT17	 EODB				1.5678E-10	2.2499E-7	14.151
LAMA3	 EODB				1.5743E-10	2.2499E-7	14.147
AREG	 EODB	5.1982	2.0151	11.482	1.6002E-10	2.2499E-7	14.131
JUNB	 EODB	3.7146	4.6077	11.469	1.6351E-10	2.2499E-7	14.111

Click on picture icon to see violin plots illustrating expression details. Click on hyperlink to see annotation details from EcoOmicsDB

Click "Proceed" to continue analysis

Previous Proceed

PLEK2

Expression

time0 time24 name

ANO MAC

time0 time24


Analysis Overview

ExpressAnalyst offers various tools to visualize results.

The analysis overview tools also provide methods to perform pathway analysis (Over-representation Analysis (ORA) and Gene Set Enrichment Analysis (GSEA))


[Home](#) > [Upload](#) > [Quality Check](#) > [Normalization](#) > [Differential Analysis](#) > [Sig. Genes](#) > Analysis Overview > [Download](#)

[Navigate to:](#)



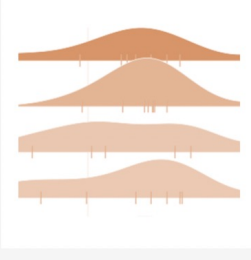
- Interactive volcano plot to display the DE genes.

[Volcano Plot](#)



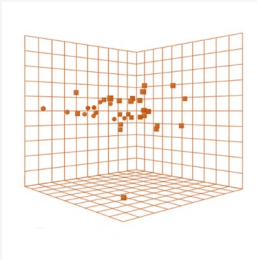
- Visualize functional categories that are enriched in a network.

[Enrichment Network](#)



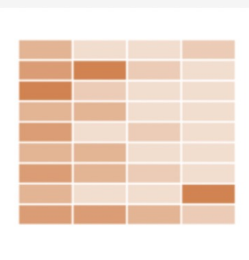
- Visualize fold-change distribution of enriched pathways

[Ridgeline Chart](#)



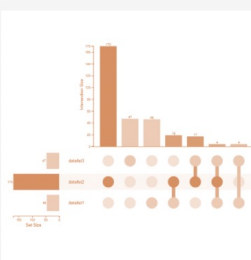
- Explore overall samples and genes in 3D space

[Dimension Reduction](#)



- Interactive heatmap to explore gene expression pattern

[ORA](#) [GSEA](#)

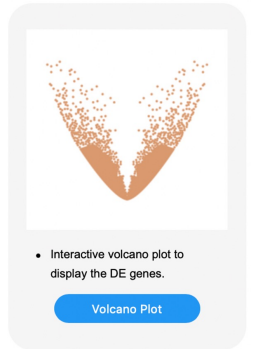


- Visualize intersections of multiple sets

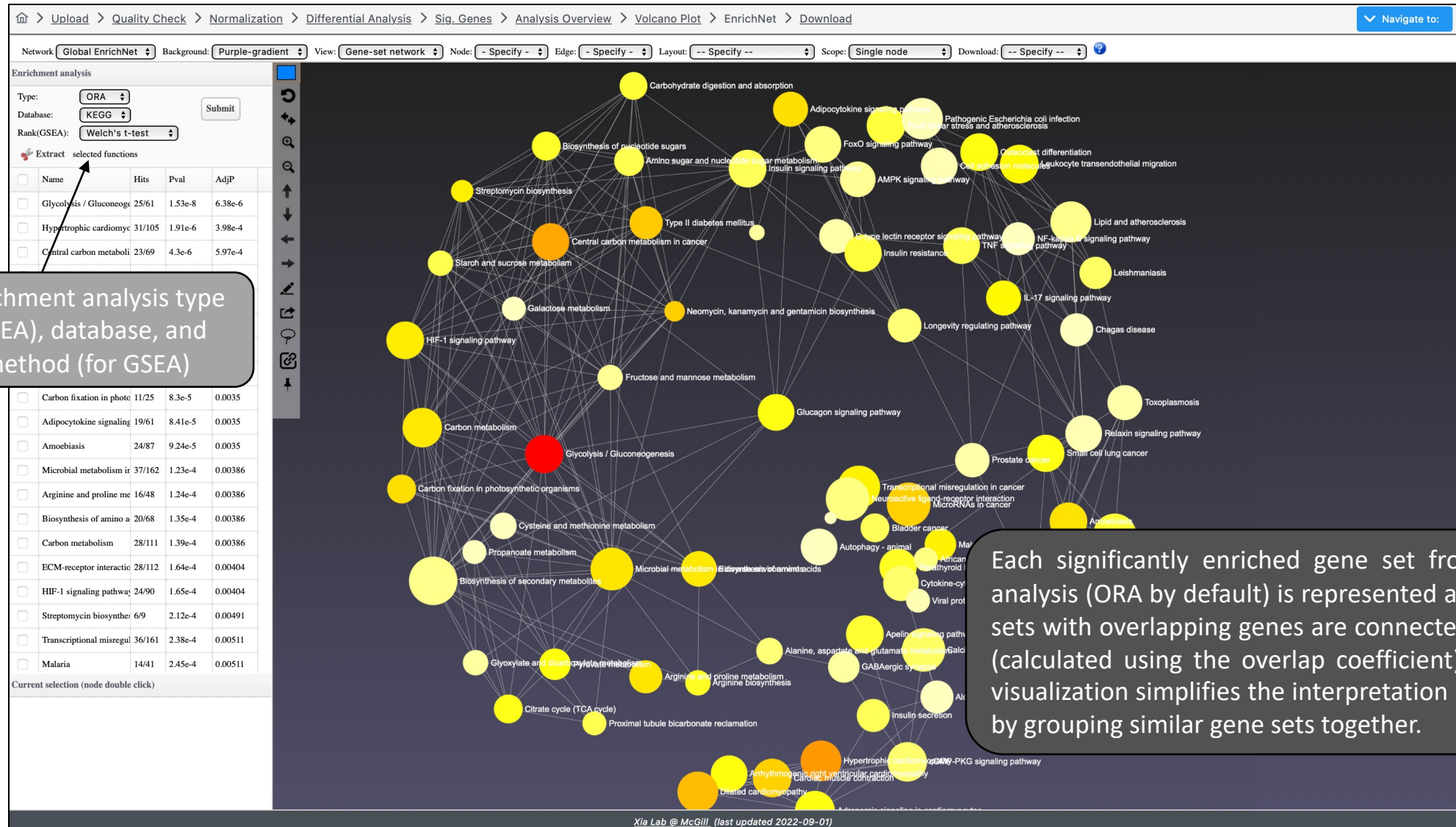
[Upset Diagram](#)

[<< Previous](#) [Downloads](#)

Analysis Overview: Volcano Plot



Analysis Overview: Enrichment Network



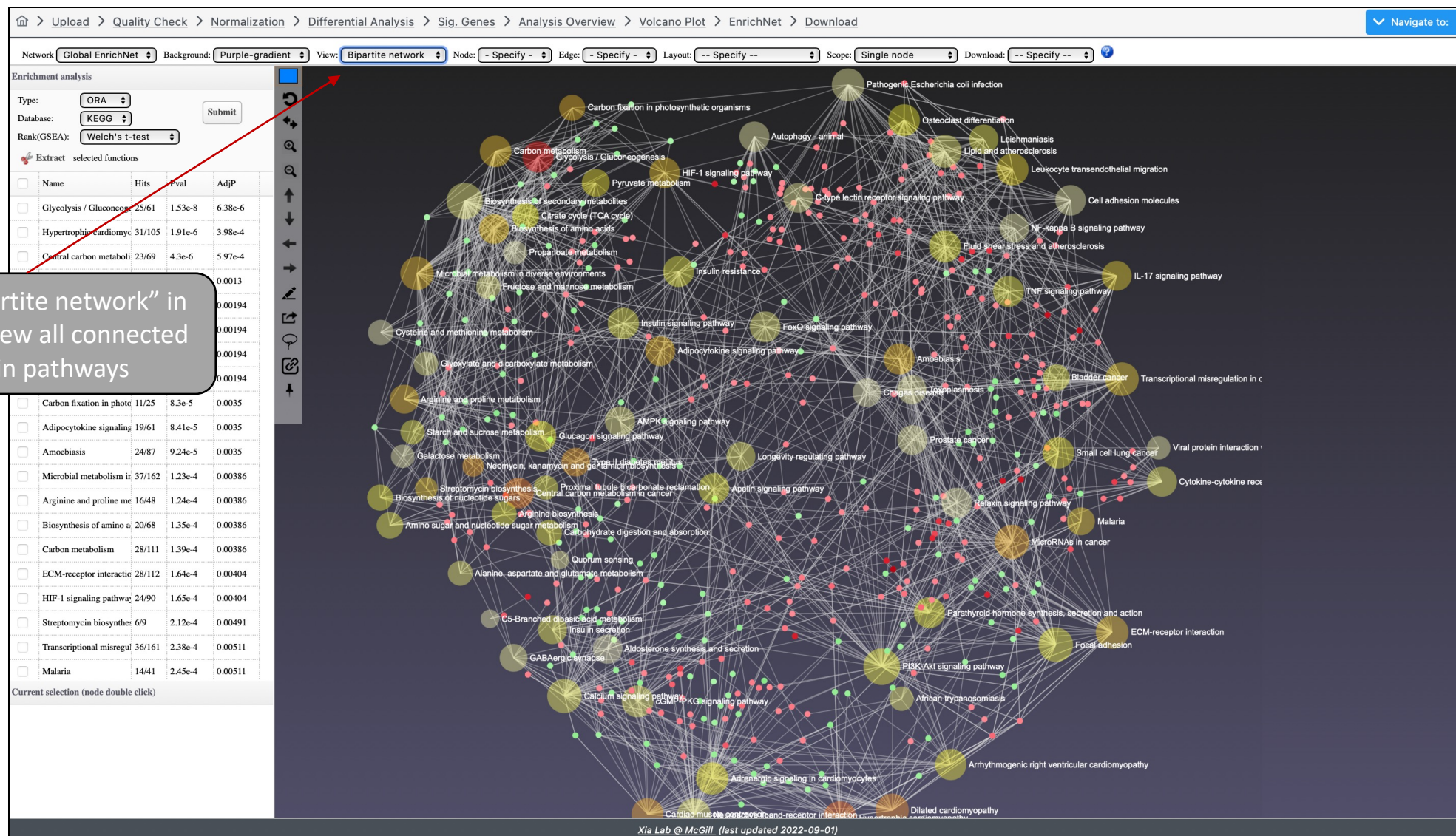
Analysis Overview: Enrichment Network



Visualize functional categories that are enriched in a network.

Enrichment Network

Select "Bipartite network" in "View" to view all connected genes in pathways

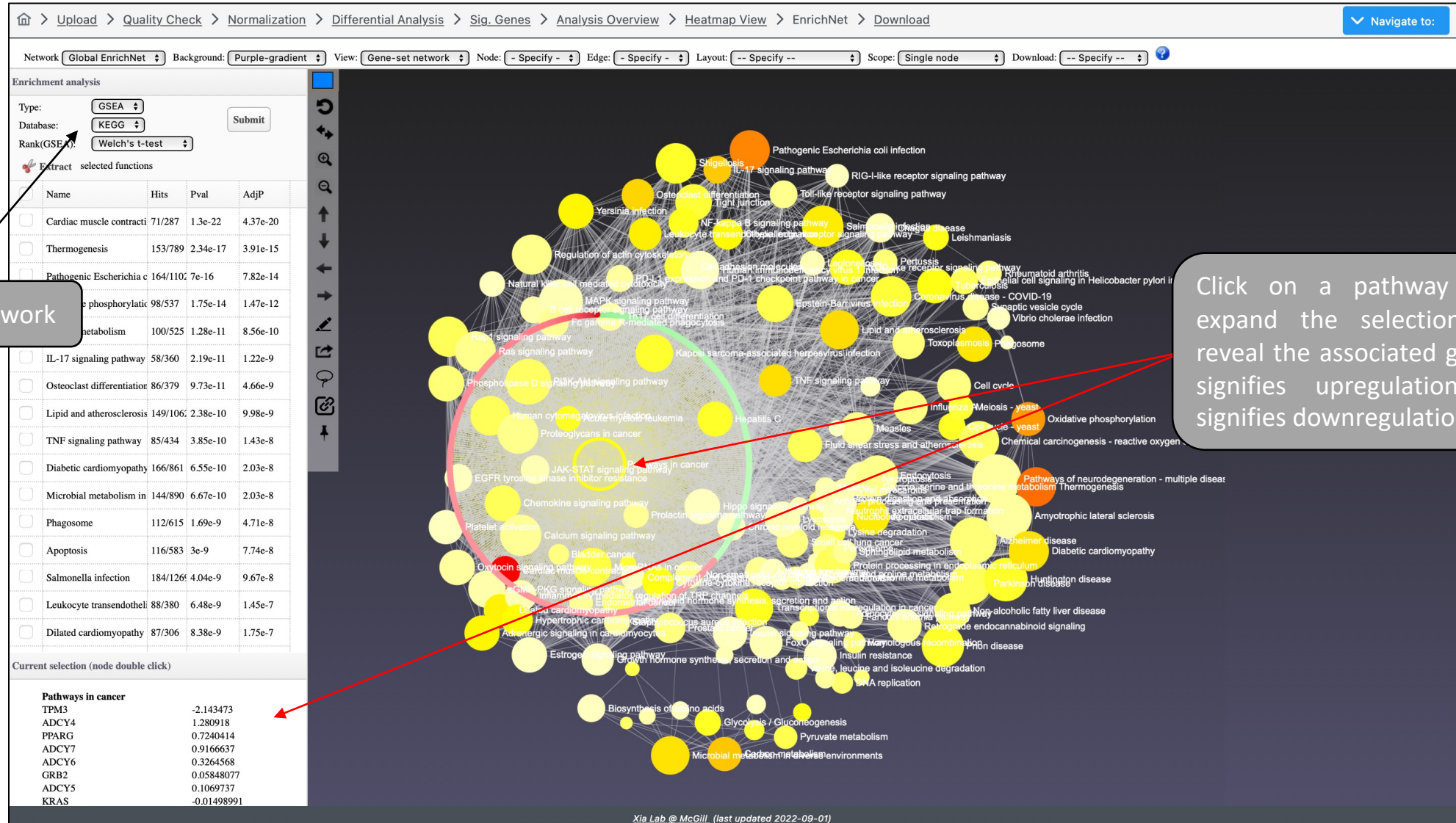


Analysis Overview: Enrichment Network

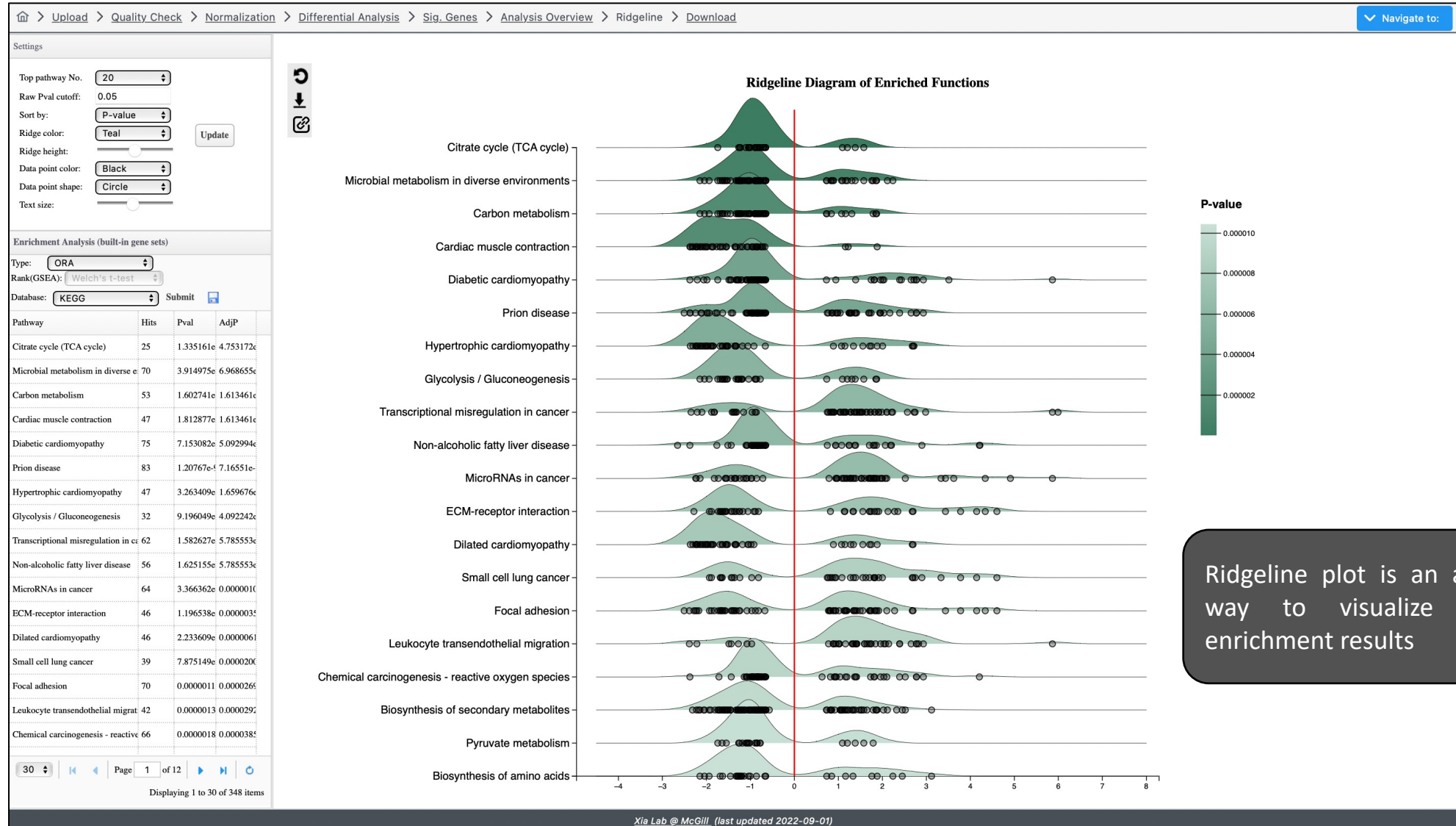
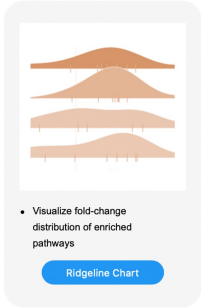


Visualize functional categories that are enriched in a network.

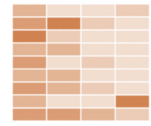
Enrichment Network



Analysis Overview: Ridgeline Plot

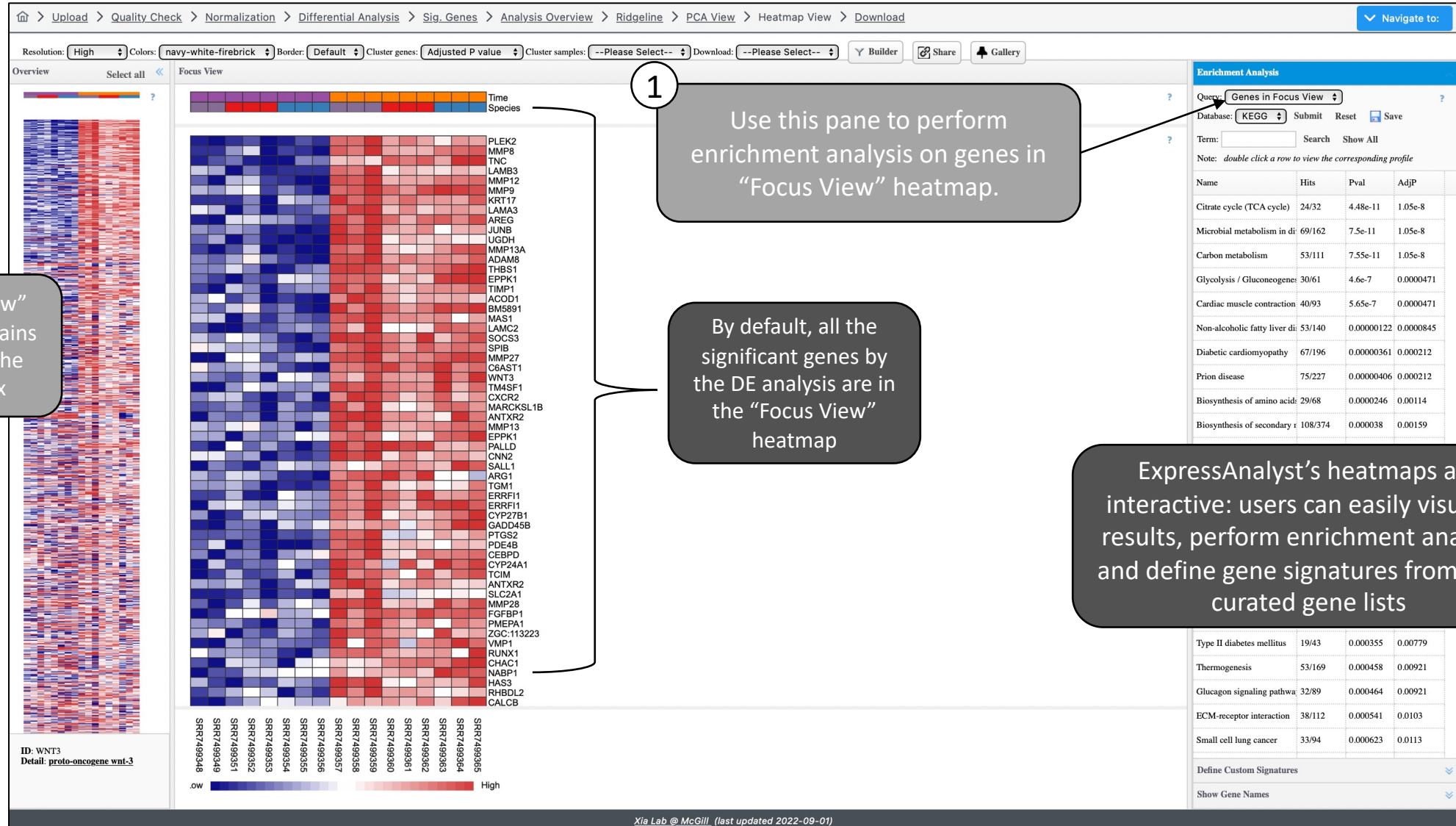


Analysis Overview: ORA Heatmap

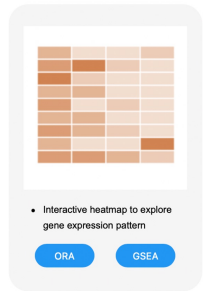


Interactive heatmap to explore gene expression pattern

ORA GSEA



Analysis Overview: ORA Heatmap



Download Results

[Home](#) > [Upload](#) > [Quality Check](#) > [Normalization](#) > [Differential Analysis](#) > [Sig. Genes](#) > [Analysis Overview](#) > [Ridgeline](#) > [PCA View](#) > [Heatmap View](#) > [Download](#)

▼ [Navigate to:](#)

Download your results

All data files and static figures that were generated during your ExpressAnalyst session are available for download below, as well as any interactive figures that you exported. Clicking on an interactive figure will navigate back to the network viewer, so you can make any last adjustments before generating figures or shareable links for publication.

[Files & Scripts](#) [Interactive Figures](#) [Static Figures](#)

Please download the results below. The **Download.zip** contains all the files in your home directory, including the static images.

Download.zip	heatmap_enrichment_2.csv
Rhistory.R	data_normalized.csv
expressanalyst_3d_pos.csv	ExpressAnalyst_pca_1.json
data_processed.csv	color_array_3.json
SigGene_time24_vs_time0_Result_1.csv	heatmap_enrichment_2.json
color_array_4.json	ridgeline_0_.csv
ridgeline_0_.json	ExpressAnalyst_matrix.json
expressanalyst_3d_load_pos.csv	counts.txt
ExpressAnalyst_loading_pca_2.json	SigGene_time24_vs_time0_Result_2.csv
ExpressAnalyst_heatmap_1.json	

Logout

Xia Lab @ McGill (last updated 2022-09-01)

The Download page compiles all figures, tables, differential expression and pathway analyses produced during the user's visit, including the raw, processed, and normalized data.

The End



For more information, visit Tutorials, Resources, and Contact pages on www.expressanalyst.ca

Also visit our forum for FAQs on www.omicsforum.ca