

# RNA-Seq Data Analysis Suite Manual

Nikhil Kumar

May 5, 2017

# Contents

|          |   |          |
|----------|---|----------|
| <b>1</b> | <b>Installation</b>                             | <b>1</b> |
| 1.1      | MongoDB . . . . .                               | 1        |
| 1.2      | Fastqc . . . . .                                | 2        |
| 1.3      | Ubuntu Packages . . . . .                       | 2        |
| 1.4      | Other software . . . . .                        | 2        |
| 1.5      | Download and Install . . . . .                  | 3        |
| 1.6      | Start the server . . . . .                      | 3        |
| 1.7      | Stop the server . . . . .                       | 3        |
| <b>2</b> | <b>Tutorial</b>                                 | <b>3</b> |
| 2.1      | Access the server . . . . .                     | 3        |
| 2.2      | Assign Job . . . . .                            | 3        |
| 2.3      | Get Significant Genes [optional] . . . . .      | 6        |
| 2.4      | Run Quality Control [required] . . . . .        | 7        |
| 2.5      | Trim Data [optional] . . . . .                  | 9        |
| 2.6      | Cluster Analysis [required] . . . . .           | 9        |
| 2.7      | Cluster Validation [optional] . . . . .         | 11       |
| 2.8      | Decision Tree [optional] . . . . .              | 12       |
| 2.9      | Manual Heatmap [optional] . . . . .             | 14       |
| 2.10     | Classify [required] . . . . .                   | 15       |
| 2.11     | Pseudotime [optional] . . . . .                 | 18       |
| 2.12     | Classified Heatmap [optional] . . . . .         | 19       |
| 2.13     | Differential Gene Analysis [optional] . . . . . | 19       |
| 2.14     | Trajectory [optional] . . . . .                 | 21       |

RNA-Seq Data Analysis Suite (RAS) is an end to end pipeline aimed to guide you through the entire RNA-Seq Analysis process. It uses many openly available bioinformatic tools.

## 1 Installation

### 1.1 MongoDB

Install mongodb from instruction on <https://docs.mongodb.com/v3.2/tutorial/install-mongodb-on-ubuntu/>

```
sudo apt-key adv --keyserver hkp://keyserver.ubuntu.com:80 --recv EA312927
echo "deb http://repo.mongodb.org/apt/ubuntu xenial/mongodb-org/3.2 multiverse"
```

Change xenial to your code name

```
sudo apt-get update
sudo apt-get install -y mongodb-org
```

Start the MongoDB service

```
sudo service mongod start
```

### 1.2 Fastqc

Download fastqc here: <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

```
chmod 755 fastqc
sudo ln -s /path/to/FastQC/fastqc /usr/local/bin/fastqc
```

### 1.3 Ubuntu Packages

```
sudo apt-get install libcurl4-openssl-dev libxml2-dev libkrb5-dev
sudo apt-get install librtmp-dev libssl-dev graphviz
```

### 1.4 Other software

1. python2
2. R
3. pip2
4. virtualenv

5. screen
6. zip
7. fastq-dump
8. fastqc
9. hisat2
10. cufflinks
11. cuffquant
12. cuffnorm
13. samtools

## 1.5 Download and Install

```
git clone https://github.com/nikhil/RAS
cd RAS
chmod +x setup.sh
./setup.sh
```

## 1.6 Start the server

After installing you need to start the server.

```
./start_server.sh
```

## 1.7 Stop the server

If you need to stop the server

```
./stop_server.sh
```

# 2 Tutorial

## 2.1 Access the server

Go to [your ip]:5000 on your server then click on 'Single-cell RNA-Seq'. Then click 'Assign Job'

## 2.2 Assign Job

Enter your name and email. Provide a name for each condition and a link for each replicate. All link and files uploaded should be a sra file. In this tutorial we will use 4 cells from GSE71485. In the GEO experiment page, you will need to click on the correct sample, click on the ftp link, find the sra file, and then copy the link. The samples in this list also include the hyperlinks for your reference.

1. Sample C10
2. Sample C110
3. Sample C111
4. Sample C118

Here is how the assign job input form should look.



## Normalized Experiment Scheduled ID: 38306



Inbox x

rutgers x



RnaSeqAnalysisSuite@rutgers.edu

9:46 PM (16 minutes ago) ☆



to me ▾

Dear Nikhil Kumar

This is a confirmation message that your experiment has been scheduled. The id for your experiment is 38306

Sincerely,  
RAS Notifier

The second email will contain links to the quality score of the samples

## Normalized Experiment Quality Check ID: 38306



Inbox x

rutgers x



RnaSeqAnalysisSuite@rutgers.edu

9:49 PM (14 minutes ago) ☆



to me ▾

Dear Nikhil Kumar

Your sequence files have been processed through quality check. You can check the results of each sample by clicking on the corresponding links below.

- SRR2132772 ( Condition: 1 Sample: 1 )  
[172.16.57.132:5000/static/data/normalize/38306/SRR2132772\\_1\\_fastqc.html](http://172.16.57.132:5000/static/data/normalize/38306/SRR2132772_1_fastqc.html)
- SRR2132692 ( Condition: 3 Sample: 1 )  
[172.16.57.132:5000/static/data/normalize/38306/SRR2132692\\_1\\_fastqc.html](http://172.16.57.132:5000/static/data/normalize/38306/SRR2132692_1_fastqc.html)
- SRR2132689 ( Condition: 4 Sample: 1 )  
[172.16.57.132:5000/static/data/normalize/38306/SRR2132689\\_1\\_fastqc.html](http://172.16.57.132:5000/static/data/normalize/38306/SRR2132689_1_fastqc.html)
- SRR2132670 ( Condition: 2 Sample: 1 )  
[172.16.57.132:5000/static/data/normalize/38306/SRR2132670\\_1\\_fastqc.html](http://172.16.57.132:5000/static/data/normalize/38306/SRR2132670_1_fastqc.html)

Sincerely,  
RAS Notifier

The third email will contain a link to the finished analysis files.

## Normalized Experiment Finished ID: 38306



Inbox x

rutgers x



RnaSeqAnalysisSuite@rutgers.edu

10:05 PM (0 minutes ago) ☆



to me ▾

Dear Nikhil Kumar

We finished processing your sequence files. You can access the output here:

[172.16.57.132:5000/static/data/normalize/38306/38306.zip](http://172.16.57.132:5000/static/data/normalize/38306/38306.zip)

Sincerely,

RAS Notifier

Note: Receiving emails from RAS has worked fine with me from Gmail. The Rutgers email seems to have some issues occasionally in receiving emails from the Sendgrid email server.

## 2.3 Get Significant Genes [optional]

If you already know which groups the cell are under you can find the genes which have the largest difference in expression across these groups. I use my experiment id that was given to me in the email from RAS.

### Experiment Significant Genes

Enter the range of numbers to signify which samples are in which group. You can use '-' (e.g. 2-10) for an inclusive range and commas to add numbers. For example 1,2,4-7 will include [1,2,4,5,6,7] for one group.

Experiment Id

Group Name

Group 1

Group Name

Group 2

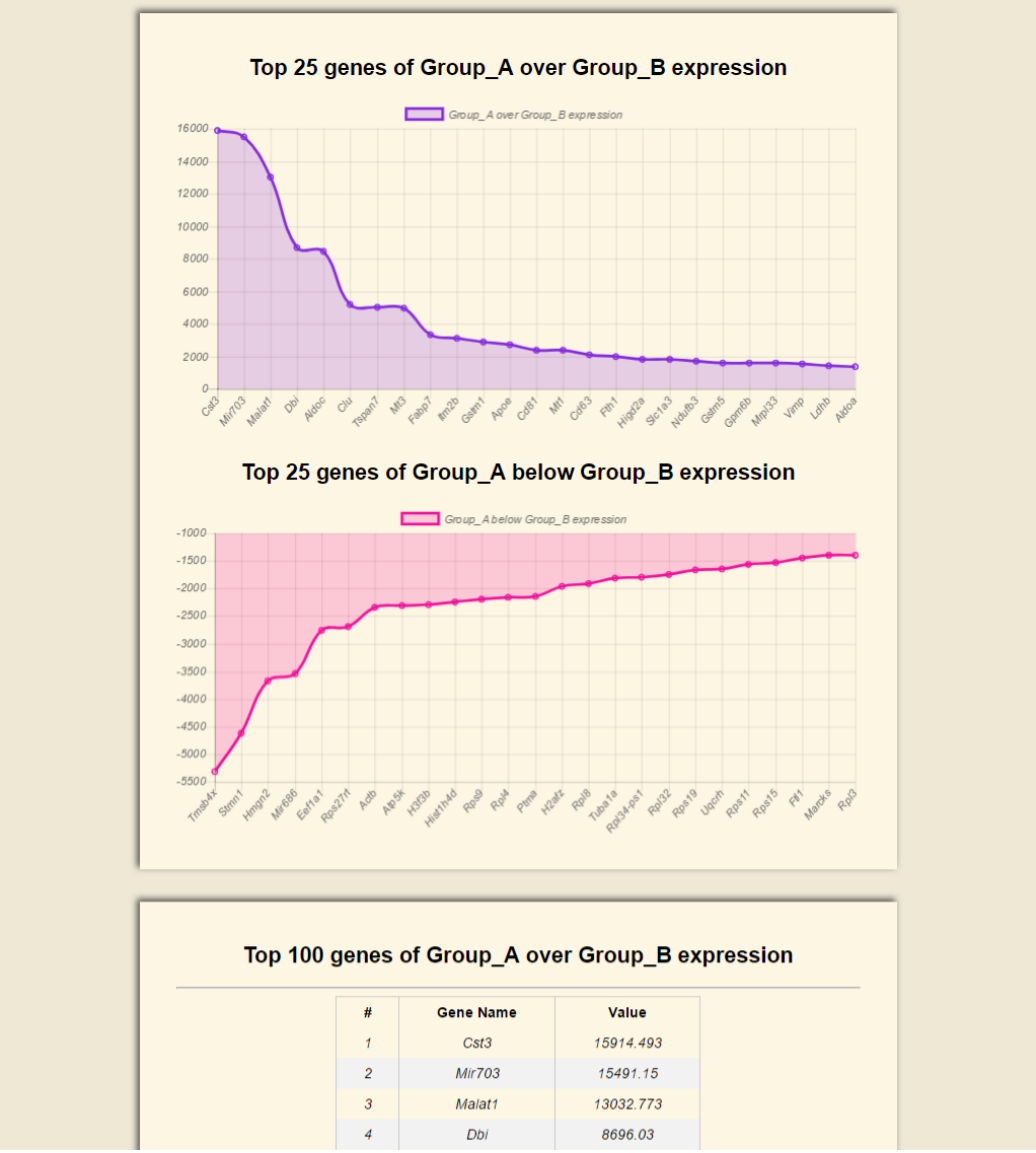
Submit

### Samples Table

| # | name          |
|---|---------------|
| 1 | Sample.C10_0  |
| 2 | Sample.C110_0 |
| 3 | Sample.C111_0 |
| 4 | Sample.C118_0 |



Here is a cropped image of the results. The results also include the top 100 genes expressed above and below the other group.



## 2.4 Run Quality Control [required]

The total mRNA in each sample is provided on the table below. In this page you need to enter the range of mRNA to use. All cells outside the range will be removed. To pass all cells use: 'inf' for the upper bound, '0' for the lower bound and '0' for minimum expression and '0' for minimum expression. Quality control works by eliminating genes that are not significantly

expressed in a certain number of cells. The gene expression and number of cells correspond to the significance and number accordingly. In the tutorial I am passing all cells. For the genes, I filter all the genes which are not expressed above the value of 1 in at least 1 of the cells.

### ☰ Quality Control

Enter the range of mRNA to use. All cells outside the range will be removed. To pass all cells use: 'inf' for the upper bound, '0' for the lower bound and '0' for minimum expression and '0' for minimum expression. Quality control works by eliminating genes that are not significantly expressed in a certain number of cells. The gene expression and number of cells correspond to the significance and number accordingly.

Experiment Id

Lower Bound

Upper Bound

Minimum Gene expression

Minimum Cells for gene expression

Submit

### Samples Table

| # | name          | Total_mRNAs      |
|---|---------------|------------------|
| 1 | Sample.C10_0  | 703422.410237183 |
| 2 | Sample.C110_0 | 617284.688602821 |
| 3 | Sample.C111_0 | 543628.101007752 |
| 4 | Sample.C118_0 | 465375.90287939  |

The confirmation page will show show the cells that passed through and the number of genes.

### ✓ Quality Control

Quality control has been completed successfully. The number of genes that passed quality control is 2725

Samples Table

| # | name          | Total_mRNAs      | num_genes_expressed |
|---|---------------|------------------|---------------------|
| 1 | Sample.C10_0  | 589596.599129317 | 2357                |
| 2 | Sample.C110_0 | 516141.766081638 | 2820                |
| 3 | Sample.C111_0 | 461941.536137503 | 2302                |
| 4 | Sample.C118_0 | 372144.055012716 | 3199                |

## 2.5 Trim Data [optional]

I will not trim data in this tutorial. However, you can use trim data to remove any outliers that were not removed from quality control. This process needs to be run after quality control.

## 2.6 Cluster Analysis [required]

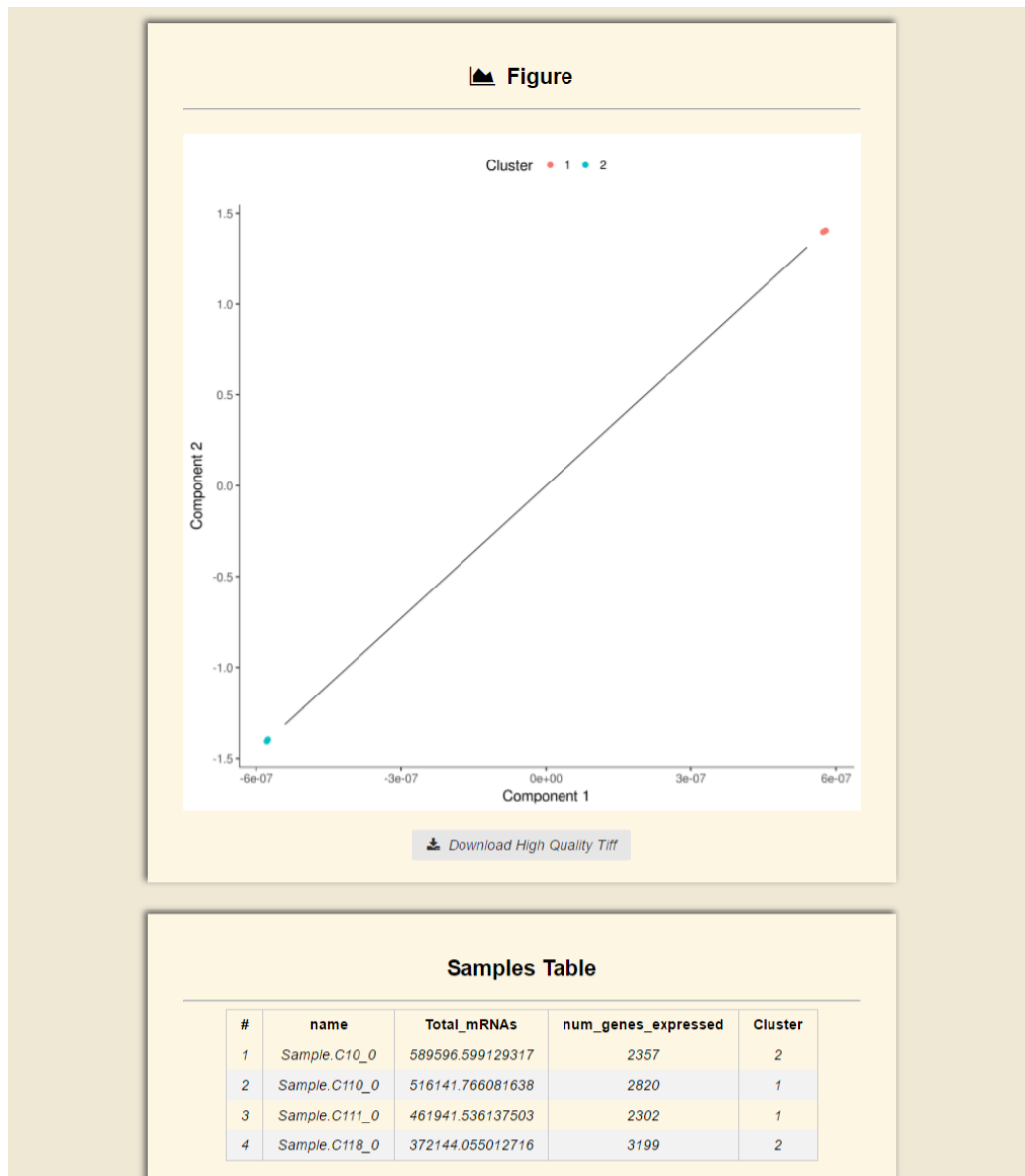
I clustered the data set into 2 groups.

### Cluster Analysis Info

---

Experiment Id

Number of Clusters



## 2.7 Cluster Validation [optional]

Often clustering analysis gives us different results at different runs. In order to confirm our results we can run the data through multiple iterations of clustering and confirm that our current clustering configuration is the most numerous one.

### ☰ Multiple Cluster Analysis Info

---

Experiment Id

User Name

User Email

Number of Clusters

Number of Iterations

| Name/Label    | Type 1 |
|---------------|--------|
| Number        | 10     |
| Sample.C10_0  | 1      |
| Sample.C110_0 | 2      |
| Sample.C111_0 | 2      |
| Sample.C118_0 | 1      |

We can see here that in all 10 iterations the cluster matched my current configuration.

## 2.8 Decision Tree [optional]

We can make a decision tree to find a significant gene which divides the group well.

### Decision Tree

Enter the range of numbers to signify which samples are in which group. You can use '-' (e.g. 2-10) for an inclusive range and commas to add numbers. For example 1,2,4-7 will include [1,2,4,5,6,7] for one group.

Experiment Id

Group Name

Group 1 Range

Group Name

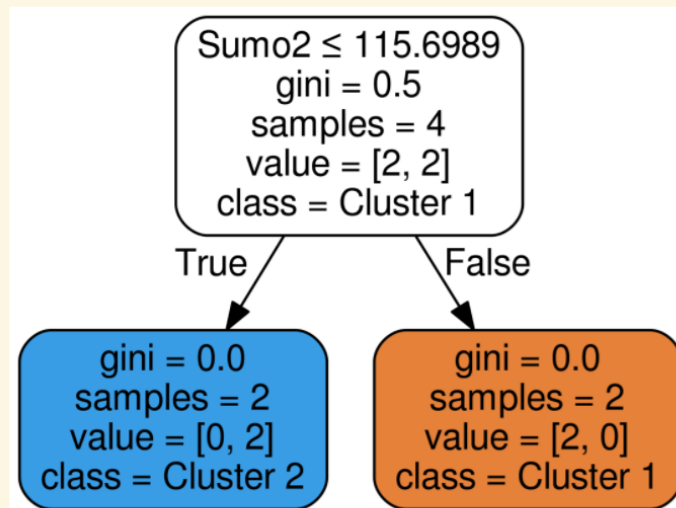
Group 2 Range

Modify Groups

### Samples Table

| # | name          | Total_mRNAs      | num_genes_expressed | Cluster |
|---|---------------|------------------|---------------------|---------|
| 1 | Sample.C10_0  | 589596.599129317 | 2357                | 2       |
| 2 | Sample.C110_0 | 516141.766081638 | 2820                | 1       |
| 3 | Sample.C111_0 | 461941.536137503 | 2302                | 1       |
| 4 | Sample.C118_0 | 372144.055012716 | 3199                | 2       |

### Figure



## 2.9 Manual Heatmap [optional]

We can get a heatmap of differential gene expression across known groups. Here is a cropped image of the results. The results also all the genes and its average expression across the groups.

### ☰ Differential Analysis Info

Enter the range of numbers to signify which samples are in which group. You can use '-' (e.g. 2-10) for an inclusive range and commas to add numbers. For example 1,2,4-7 will include [1,2,4,5,6,7] for one group. The min mean difference removes the genes where the mean gene expression difference from any two classes is less than the specified value.

Experiment Id

Number of Genes

Min Mean Difference

Group Name

Group 1 Range

Group Name

Group 2 Range

Modify Groups

### Samples Table

| # | name          | Total_mRNAs      | num_genes_expressed | Cluster |
|---|---------------|------------------|---------------------|---------|
| 1 | Sample.C10_0  | 589596.599129317 | 2357                | 2       |
| 2 | Sample.C110_0 | 516141.766081638 | 2820                | 1       |
| 3 | Sample.C111_0 | 461941.536137503 | 2302                | 1       |
| 4 | Sample.C118_0 | 372144.055012716 | 3199                | 2       |





## 2.10 Classify [required]

I already know that the cells in cluster 2 are active neural stem cells. I can use the critical gene I found from the decision tree to classify the dataset through gene expression.

User Info

Experiment Id38306

Label 1

Label NameActive

Condition 1Sumo2<=115.69

Modify conditionsAdd conditionRemove condition

Label 2

Label NameQuiescent

Condition 1Sumo2>115.69

Modify conditionsAdd conditionRemove condition

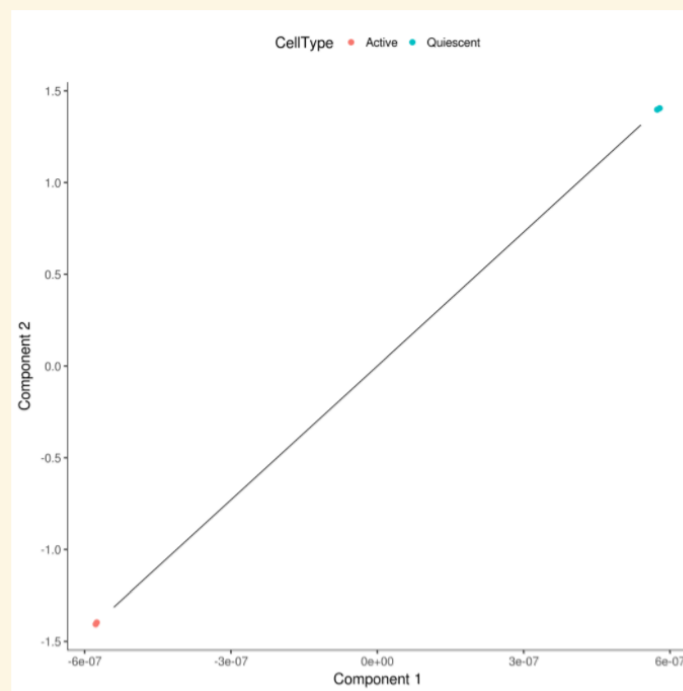
Modify labelsAdd labelRemove label

SubmitSubmit

### ✓ Classification

Classification has completed successfully.

### Figure



[Download High Quality Tiff](#)

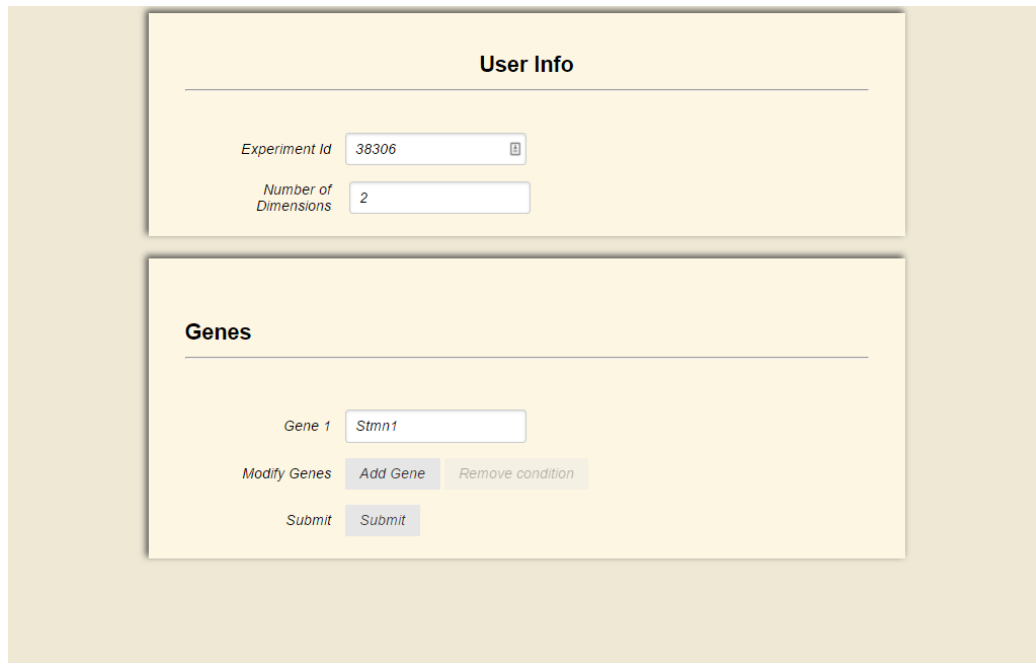
### Samples Table

| # | name          | Total_mRNAs      | num_genes_expressed | Cluster | CellType  |
|---|---------------|------------------|---------------------|---------|-----------|
| 1 | Sample.C10_0  | 589596.599129317 | 2357                | 2       | Active    |
| 2 | Sample.C110_0 | 516141.766081638 | 2820                | 1       | Quiescent |
| 3 | Sample.C111_0 | 461941.536137503 | 2302                | 1       | Quiescent |
| 4 | Sample.C118_0 | 372144.055012716 | 3199                | 2       | Active    |

If you choose not to classify through gene expression you will need to classify manually.

## 2.11 Pseudotime [optional]

You can compare gene expression with the cells ordered in pseudotime. For more information on pseudotime ordering, it is recommended to read the Monocle Manual.

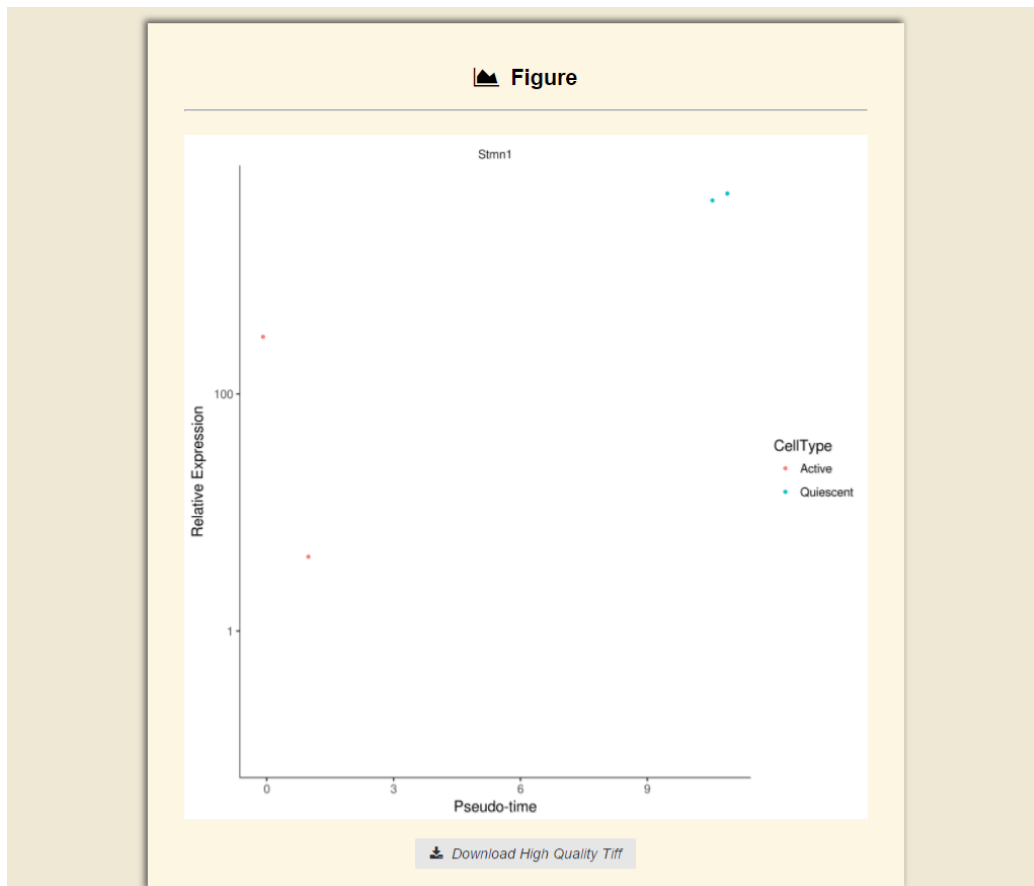


The screenshot displays the Monocle web interface with two main sections: 'User Info' and 'Genes'. The 'User Info' section contains input fields for 'Experiment Id' (38306) and 'Number of Dimensions' (2). The 'Genes' section includes a 'Gene 1' input field with 'Stmn1', a 'Modify Genes' section with 'Add Gene' and 'Remove condition' buttons, and a 'Submit' button.

| User Info            |       |
|----------------------|-------|
| Experiment Id        | 38306 |
| Number of Dimensions | 2     |

| Genes        |   |
|--------------|---|
| Gene 1       | Stmn1   |
| Modify Genes | <button>Add Gene</button> <button>Remove condition</button> |
| Submit       | <button>Submit</button>                                     |



## 2.12 Classified Heatmap [optional]

This is the same as the manual heatmap. The only difference is that the group and group name will be set based on the classification.

## 2.13 Differential Gene Analysis [optional]

This shows the gene expression of a specific gene across all the classified groups.

User Info

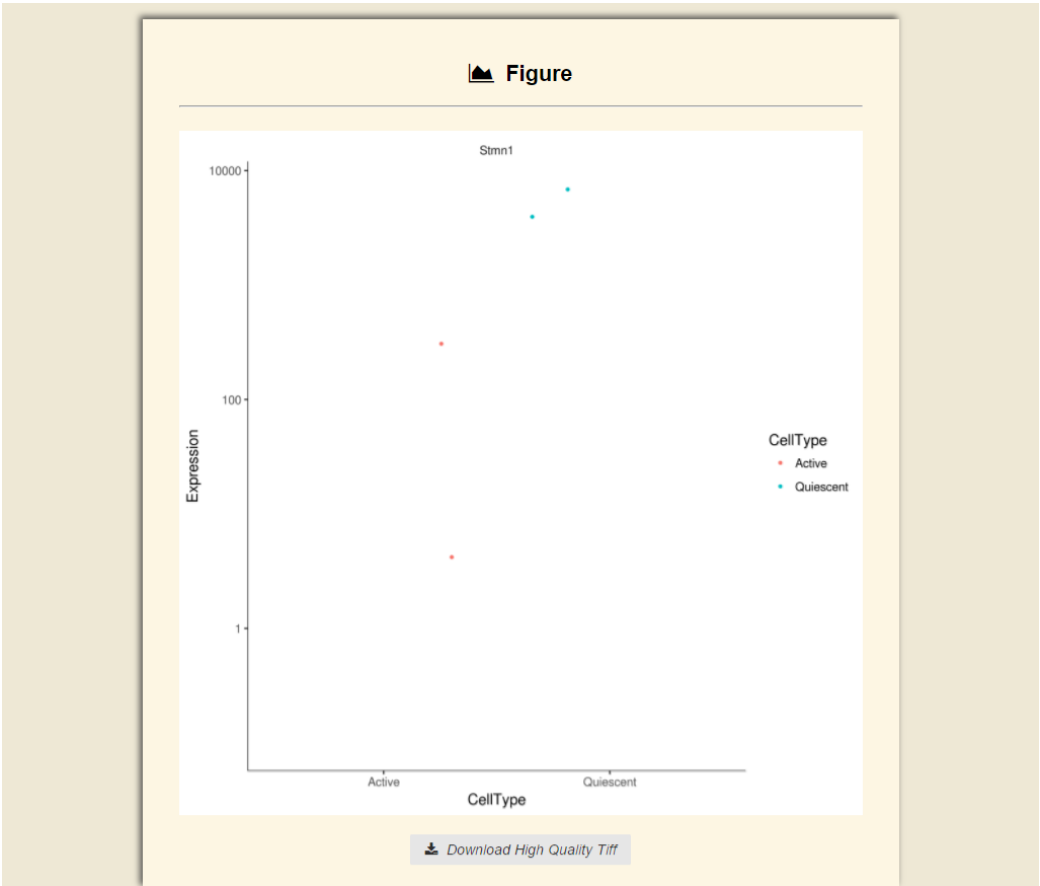
Experiment Id38306

Genes

Gene 1Stmn1

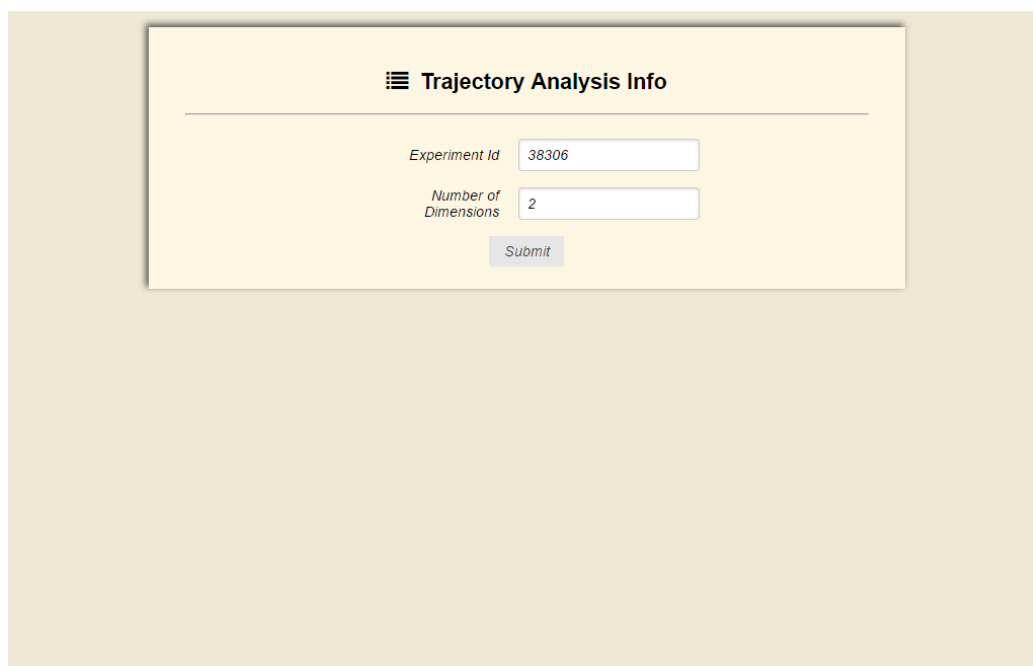
Modify GenesAdd GeneRemove condition

SubmitSubmit




## 2.14 Trajectory [optional]

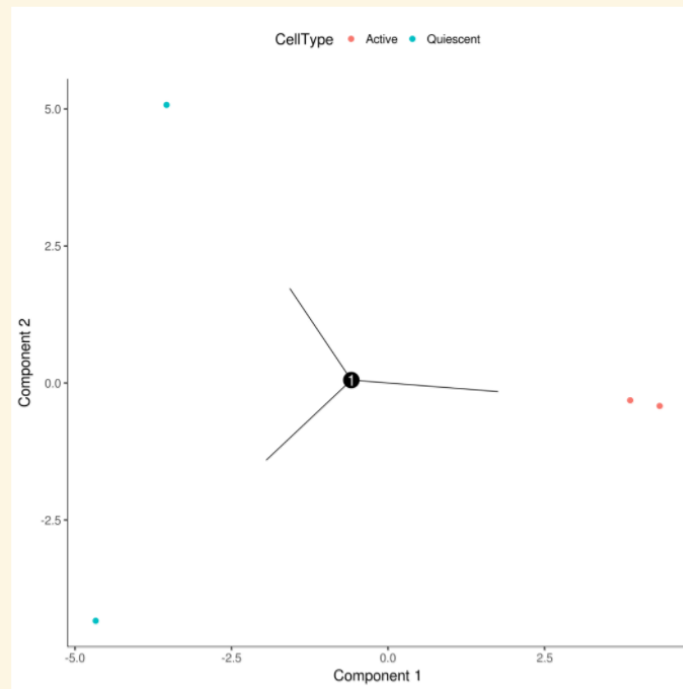
The cells can be ordered in a trajectory. For more information refer to the Monocle manual.




The screenshot shows a web interface for "Trajectory Analysis Info". It features a title bar with a hamburger menu icon and the text "Trajectory Analysis Info". Below the title bar, there are two input fields: "Experiment Id" with the value "38306" and "Number of Dimensions" with the value "2". A "Submit" button is located below the input fields. The entire form is set against a light beige background.

| Trajectory Analysis Info              |       |
|---------------------------------------|-------|
| Experiment Id                         | 38306 |
| Number of Dimensions                  | 2     |
| <input type="button" value="Submit"/> |       |

 **Figure**



 [Download High Quality Tiff](#)