

USER GUIDE

INSERM 33

Yohan Lefol

Alexis Dupis

Ugo Vidal

Marie Terrien

January 13, 2020



1 Introduction

The non-canonical analysis tool (INSERM 33) is a tool with several uses. It's main feature is a non-canonical gene analysis from a gene list. Yet this tool also boasts several other interesting functionalities such a differential expression analysis, customizable MA and Volcano plots which in turn allow to obtain significant gene sets according to a users custom parameters. This guide will go over each functionality and explain how to best use **THE TOOL**.

As can be seen in Figure 1, the application is split into several tabs. The information tab shows two pdf files, the first being this one, and the second being a detailed explanation of the tool and it's components. The last tab 'Citation' is a simple tab which gives the citation information for the use of this tool. This guide will now go over the remaining tabs and explain their uses and how to maneuver them.



Figure 1: The main tab for the application.

2 Installation

HOW DOES THE INSTALLATION WORK.

3 Connection to database

The identification of non-canonical genes is the main aspect of this tool. As the list of non-canonical genes is increasing as discoveries are being made, we could not 'hard code'

a list of genes in the program as we wanted to ensure that the search list is being updated without having users re-download the application for each update. Therefore we implemented a database that can be updated by the owners of the tool (Cellomet) and that can be accessed remotely through this application. In order to establish this connection, a user needs to fill in the small questionnaire as seen in Figure 2, and click the 'Connect to database' button. We created the small questionnaire in order to see the main users of this tool as we intend to keep developing new functionalities and we want these functionalities to be catered to the types of study that require it the most.

The screenshot shows a web browser window with the address bar displaying '127.0.0.1:7636'. The application title is 'INSERM 33'. On the left is a dark sidebar with a menu containing: Information, Connect to database (active), Run analysis with DESeq2, Analyze a gene list, Results, Custom MA plots, Custom Volcano plots, and Citation. At the bottom of the sidebar is an 'Exit App' button. The main content area is titled 'Establish connection to the database' and contains the following form elements: an 'Email Address' text input field, a 'Type of Study' dropdown menu with 'Metabolic' selected, a 'Position held' dropdown menu with 'Researcher' selected, and a 'Comments' section with a text area containing the placeholder 'Feel free to add a comment'. At the bottom of the form is a 'Connect to database' button.

Figure 2: The connect to database tab of the application.

The database connection is only necessary if a non-canonical gene analysis is to be performed, either individually or as part of a DESeq2 analysis. If this is not the users intentions, he does not need to establish a connection to the database.

4 DESeq2 analysis

DESeq2 is a powerfull tool which allows users to perform differential analyses using RNAseq 'COUNTS' data (Love, Huber, & Anders, 2014). However this tool requires a very specific data format. This part of the guide will go over the data format necessary to run a DESeq2 analysis, the added parameters specific to this application, and the results that are generated.

4.1 Tab layout

As seen in Figure 3, there are several requirements for the the DESeq2 analysis. First we observe that two csv files are needed, these will be the two files representing the two conditions that will be analyzed. We then observe a small check box which asks if the user want to perform a non-canonical gene analysis. If this box remains checked, the user will have to be connected to the database and the non-canonical gene analysis will be run on the list of significant genes isolated by the DESeq2 analysis (explained below). A user can also select a directory, this will be the folder in which the several DESeq2 results will be saved. There are then two text boxes which can be customized to be whatever text a user wishes, the content of these text boxes is used to name the files that will be produced by the analysis. Condition 1 represents File 1 and Condition 2 represents File 2. Lastly there are two check boxes that check if the user would like to create the standard MA and Volcano plots. If checked, two types of each plot will be made, one type with text, and the other without text. The plots generated will be with standard values, however a user will be able to take the differential expression data obtained from the DESeq2 analysis and make his own customized MA and Volcano plots with the custom plot tabs, see section 7 for more details.

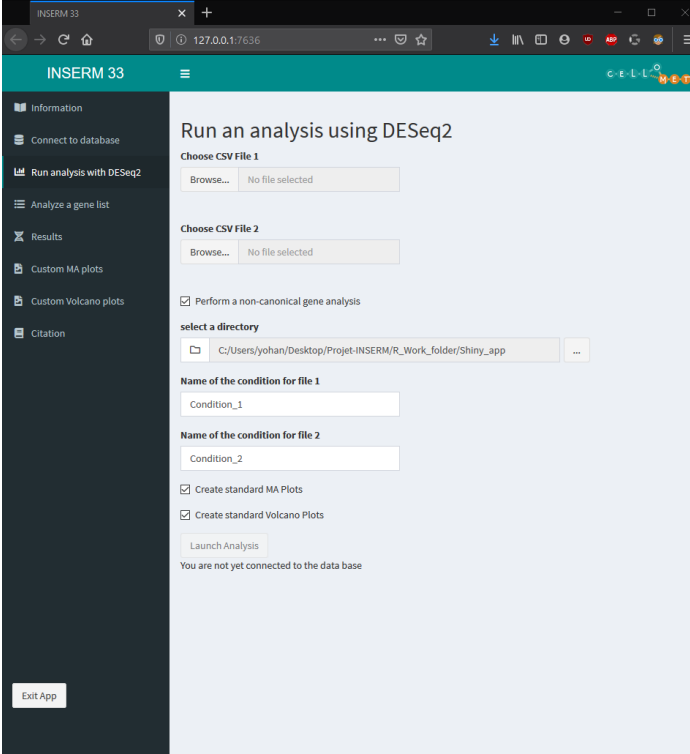
The image shows a web browser window with the title 'INSERM 33'. The address bar shows '127.0.0.1:7636'. The page has a dark teal header with the text 'INSERM 33' and a hamburger menu icon. On the left is a dark sidebar with a menu: 'Information', 'Connect to database', 'Run analysis with DESeq2' (highlighted), 'Analyze a gene list', 'Results', 'Custom MA plots', 'Custom Volcano plots', and 'Citation'. At the bottom of the sidebar is an 'Exit App' button. The main content area is light blue and titled 'Run an analysis using DESeq2'. It contains two 'Choose CSV File' sections, each with a 'Browse...' button and 'No file selected' text. Below these is a checked checkbox 'Perform a non-canonical gene analysis'. Then, a 'select a directory' section with a file path 'C:/Users/yohan/Desktop/Projet-INSERM/R_Work_folder/Shiny_app' and a folder icon. Two text input fields are labeled 'Name of the condition for file 1' (with 'Condition_1' entered) and 'Name of the condition for file 2' (with 'Condition_2' entered). At the bottom are two checked checkboxes: 'Create standard MA Plots' and 'Create standard Volcano Plots', followed by a 'Launch Analysis' button. A message at the very bottom states 'You are not yet connected to the data base'.

Figure 3: The main screen for the DESeq2 analysis tab.

4.2 Data format

As previously stated, DESeq2 requires a rather specific data format. We will be using Figure 4 to explain the format. This format is presented as a CSV file, meaning a comma separated

vector, or more simply put, a file in which the different elements/values are separated by a comma. In order to better explain the data format and what a csv file really is, an excel representation of Figure 4 has been created and is shown in Figure 5

```
Gene,01daeea0-2a6d-450e-b541-5e403e998637,05a38163-a505-4f97-b2a3-7430c49b14f4,1
TSPAN6,4280,5096,1949,2711,1669,3992,5499,5770,794,4039,4943,4236,21
TNMD,4,8,2,3,2,7,14,5,2,26,17,58,4,8,3,1933,2,27,31,2,10,1
DPM1,1293,1019,440,1088,844,1271,1137,1546,440,1549,1319,1374,753,
SCYL3,728,563,127,458,650,542,651,742,161,596,811,453,547,551,519
C1orf112,565,342,123,223,480,517,500,745,151,478,477,184,270,378,:
```

Figure 4: The data format for a DESeq2 analysis.

When comparing both Figure 4 and Figure 5 we can clearly observe that the commas create the 'separation' of the values. Now we will explain what these values are. The first column represents the genes, it is important to note that these genes are in gene symbols and not another type of format such as Ensembl. DESeq2 will perform the analysis regardless of the type of gene name used however **a non-canonical gene analysis can only be performed with gene symbols**. As such it is advised to use gene symbols in the data in order to fully benefit from this tool. Next we have the other columns which represents the different sequencings. In this example, the second column shows the RNAseq 'COUNTS' results of patient #1 and the second column shows the same but for patient #2.

	A	B	C
1	Gene	01daeea0-2a6d-450e-b541-5e403e998637	05a38163-a505-4f97-b2a3-7430c49b14f4
2	TSPAN6	4280	5096
3	TNMD	4	8
4	DPM1	1293	1019
5	SCYL3	728	563
6	C1orf112	565	342
7	FGR	762	1027
8	CFH	480	1037

Figure 5: A simplified view of the DESeq2 data format.

That's it, that is the data format required. As previously mentioned, two files are necessary, i.e: one file per condition. An easy example to explain this would be the comparative study of young vs old pulmonary cancer patients. In order to do such an analysis, we would have file #1 contain all of the patients in the 'young' category and file #2 would contain all the patients of the 'old' category. Of course both files must respect the data format.

4.3 launching the analysis

Once all necessary files have been uploaded, a user can click the 'launch analysis' button, at which point a pop-up will appear indicating the the analysis is being performed. This pop-up will lock the application preventing any further use, this is done to avoid the application from crashing if too many actions are performed at once. It is important to note that a DESeq2

analysis can be very quick just as it may take several hours. The length of the analysis depends on the size of the files used as well as the computing strength of the users computer. If a user is using a old computer, it may be advised to let the computer run the analysis without doing anything else on the side. If an error occurs during the analysis, another pop-up will appear stating that there was an error and that the user can read about the error in the error_log.txt that was created. Using the information from the error log, the user should consult section 8 in order to solve the issue. Once the analysis is done and if no error has occurred, the pop-up message will be replaced by a different pop-up message indicating that the analysis is done, the pop-up message will also remind the user in which files he has stored the results of the analysis.

4.4 Results obtained

Several results can be obtained from this analysis, some are obtained regardless of user input, others will only appear if a user has asked for those results to be produced.

4.4.1 Main results

A DESeq2 analysis generates three standard results:

- `diffexpr_results condition_1 vs condition_1 .csv`
A differential expression file
- `condition_1vscondition_1_RESULTS_VOLCANO.csv`
A differential expression file for significant genes
- `gene_list_ condition_1 vs condition_1 _Most_Significant.txt`
A significant gene list

It is important to note that the significant genes were obtained by filtering the main differential expression file for standard significance parameters while using a volcano plot to visualize the significant genes. The standard significance algorithm is the following:

$$padj < 0.05 \ \& \ abs(log2FoldChange) > 2$$

This equation indicates that genes which have a Padj value below 0.05 and a log2FoldChange above 2 or below -2 are considered to be significant. A user can also generate his own list of significant genes by modifying these parameters. This is further explained in section 7 of this guide.

4.4.2 Figures generated

If a user has selected that he wanted standard MA and Volcano plots, there will be four png files in the results folder; Two MA plots, one with text and the other without text, the same applies for the volcano plot. The plots with texts are much larger in order to accommodate for the additional text besides the points representing significant genes. If a user wishes to modify these plots and create his own set, it can be easily be done in the custom plot tabs, further explained in section 7 of this guide.

4.4.3 Non-canonical gene analysis

If a user had selected that he wanted a non-canonical gene analysis to be performed, there will be three extra text files in the results folder. That is, if any non-canonical genes were found within the significant genes list. The files are the following:

- `non_canonic_results.txt`
This file will contain the gene symbols and gene names that were found to be non-canonical genes within the significant gene list that was analyzed. The file also contains the non-canonical action and localization of the genes identified
- `canonic_results.txt`
Similarly to non-canonical results, however this one contains the canonical action and localization.
- `references.txt`
Any and all references that were used to identify non-canonical genes and their actions.

The results were saved as text files that can easily be read by Microsoft Excel, or any similar program. Additionally, results can be viewed immediately in the results tab of the application, see section 6 of this guide for more information.

5 Gene list analysis

This tab allows users to perform a non-canonical gene analysis with a simple gene list file. The only obligation that this file must respect is that the genes must be in a column with a column header named 'Genes'. If the file entered is correct, a table will appear showing that specific column as seen in Figure 6. If this does not occur, a text will appear, stating that a different file should be chosen.

Considering that the requirement for this analysis is a single column carrying the name of 'Genes', possible files can be the significant gene list obtained from the DESeq2 analysis, or the significant differential expression results obtained from either the DESeq2 analysis or the downloading of the data from custom plots in section 7.

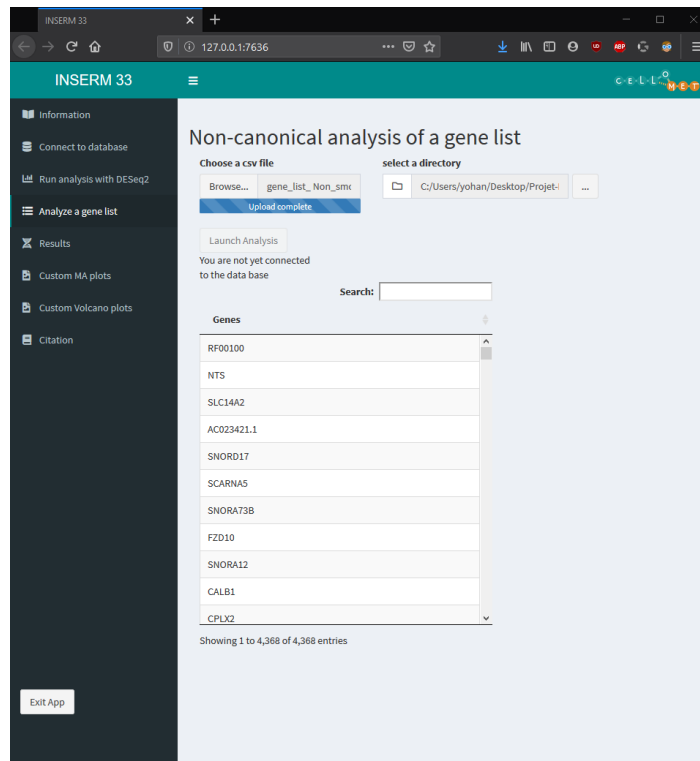


Figure 6: The gene list analysis tab once a correct file has been uploaded.

6 Results

The results tab, seen in Figure 7, shows the results of a non-canonical gene analysis within the application itself. There isn't much to be said about this tab, it is only for visualization. Every registered non-canonical genes are registered in a database hosted on Cellomet.com which can be accessed, updated and generally modified by the owner of Cellomet. For any questions regarding the non-canonical genes, their function, localization etc. Please contact Cellomet via their website.

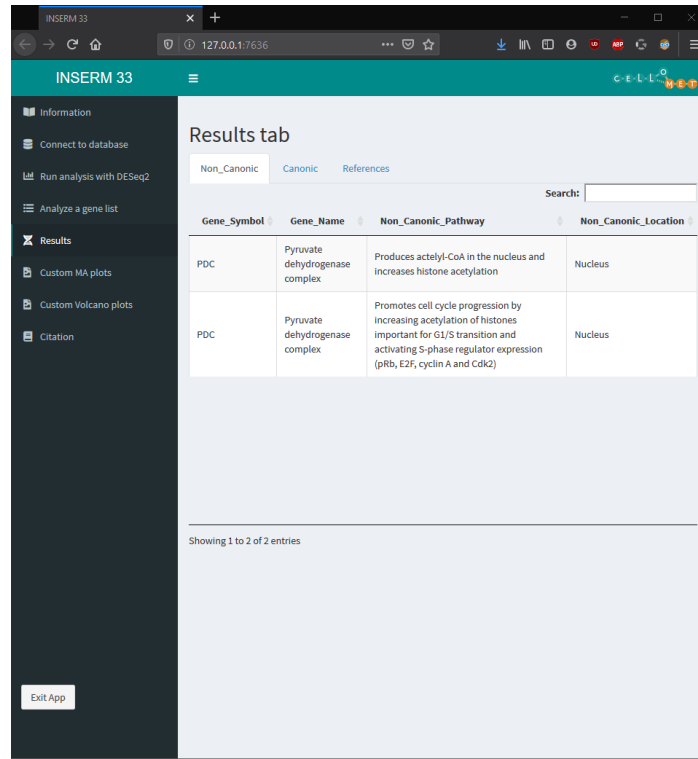


Figure 7: The results tab after a non canonical gene list analysis.

7 Custom MA and Volcano plots

This application allows users to create their own custom figures, and from these figures, significant gene data can be downloaded. There are two possible figures, MA plots and Volcano plots, each plot type has its own tab, however they behave very similarly, the volcano plot tab can be seen in Figure 8. The main thing to know about these tabs is that a differential expression file is required, like the one obtained from DESeq2 analyses. The columns required are the 'Genes' column, the 'baseMean' column, the 'log2FoldChange' column, the 'pvalue' column, and the 'padj' column, each written the same way as they are written in the guide. The file must also be a csv file. Aside that, each element of the different tabs is self explanatory. After each parameter modification, the create plot button needs to be clicked again. When a plot is loading, the rest of the application is locked to prevent the application from crashing due to a potential large number of buffered actions. On every button click the plot must be re loaded from scratch, this factor should be taken into consideration if a user is using a large differential expression files.

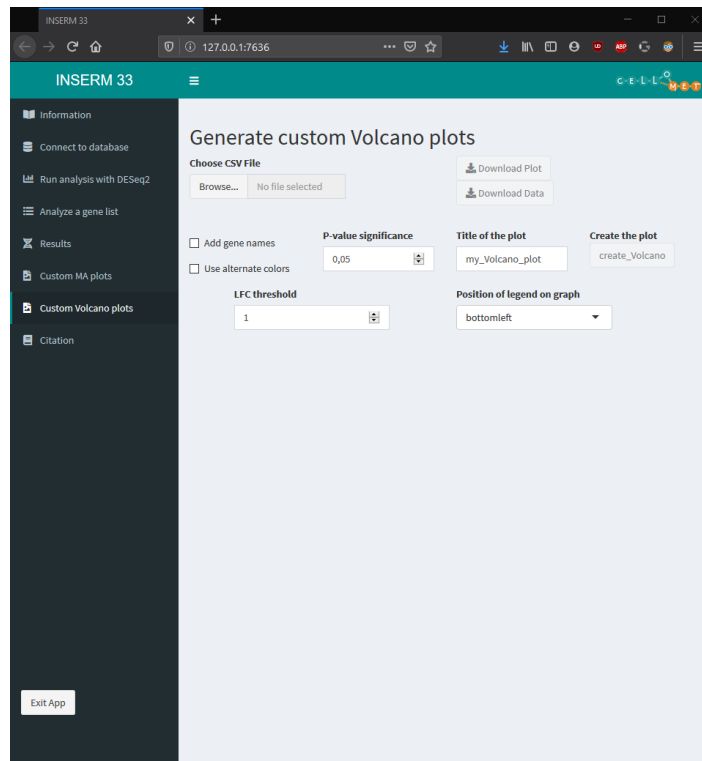


Figure 8: The tab to create a custom volcano plot.

8 Common errors

References

Love, M. I., Huber, W., & Anders, S. (2014). Moderated estimation of fold change and dispersion for rna-seq data with deseq2. *Genome biology*, 15(12), 550.