

User's Guide for OneStopRNAseq

Author List and Affiliations

Rui Li, Department of Molecular, Cell and Cancer Biology, University of Massachusetts Medical School, 364 Plantation Street, Worcester, MA 01605, USA. Rui.Li@umassmed.edu

Kai Hu, Department of Molecular, Cell and Cancer Biology, University of Massachusetts Medical School, 364 Plantation Street, Worcester, MA 01605, USA. Kai.Hu@umassmed.edu

Haibo Liu, Department of Molecular, Cell and Cancer Biology, University of Massachusetts Medical School, 364 Plantation Street, Worcester, MA 01605, USA. Haibo.Liu@umassmed.edu

Michael R. Green, Department of Molecular, Cell and Cancer Biology, University of Massachusetts Medical School, 364 Plantation Street, Worcester, MA 01605, USA. michael.green@umassmed.edu

Corresponding author: Lihua Julie Zhu, Department of Molecular, Cell and Cancer Biology, Department of Molecular Medicine, Program in Bioinformatics and Integrative Biology, University of Massachusetts Medical School, 364 Plantation Street, Worcester, MA 01655, USA. Julie.Zhu@umassmed.edu

RL, KH, and HL contributed equally to this work.

Overview

OneStopRNAseq (<https://mccb.umassmed.edu/OneStopRNAseq>) is an easy to use web application designed for comprehensive analyses of RNA-seq data for biologists and bioinformaticians. In order to simplify and streamline RNA-seq analysis, we integrated many widely used analysis components into our pipeline, including differential gene expression (DGE) analysis, differential exon usage (DEU) analysis, differential alternative splicing (DAS) analysis, gene set enrichment analysis (GSEA), differential transposon element (DTE) analysis, and allele-specific expression (ASE) quantification. To be flexible and shorten analysis time, OneStopRNAseq provides multiple analysis entry points, i.e., users can start their analyses with raw FASTQ files, aligned BAM files, gene count table files or rank files according to the analysis goals and available files as shown in the right panel of Figure 1. In summary, users can simply put their data in Dropbox, provide Dropbox download links and sample information (metadata)

through the website, select analyses of interest, and comparisons/contrasts of interest, submit the job, and get analysis results. Data uploading is carried out in batch via Wget and the users do not have to upload one file at a time via the website to avoid any interruption by unstable internet connection and long wait time for the users. Alternatively, users can input GEO accession number and the associated metadata will be automatically retrieved from the database. After users review and verify associated metadata, raw data will be downloaded and extracted using prefetch and fastq-dump automatically. Therefore, OneStopRNAseq is developed not only for analyzing your own data but also convenient for analyzing public datasets.

In the following section, we provide more details on how to use our web application. We will keep updating the user's guide with new functionalities added. For the most recent version, we encourage you to visit our web site.

Web interface

Figure 1 shows the homepage of our web application. We listed the currently supported analyses on the left panel where you can hover over each item to view the software packages that are being used. The main workflow is illustrated on the right panel where all of the black boxes (and round boxes) are “clickable”. The clickable boxes are possible choices of entry points to the pipeline. To start with “FASTQ” or “BAM” files, you have to register yourself by filling in your email and other information in the registration form, whereas no registration is needed to start the analysis from “Count table”, “RNK file”, or “Visualization”.

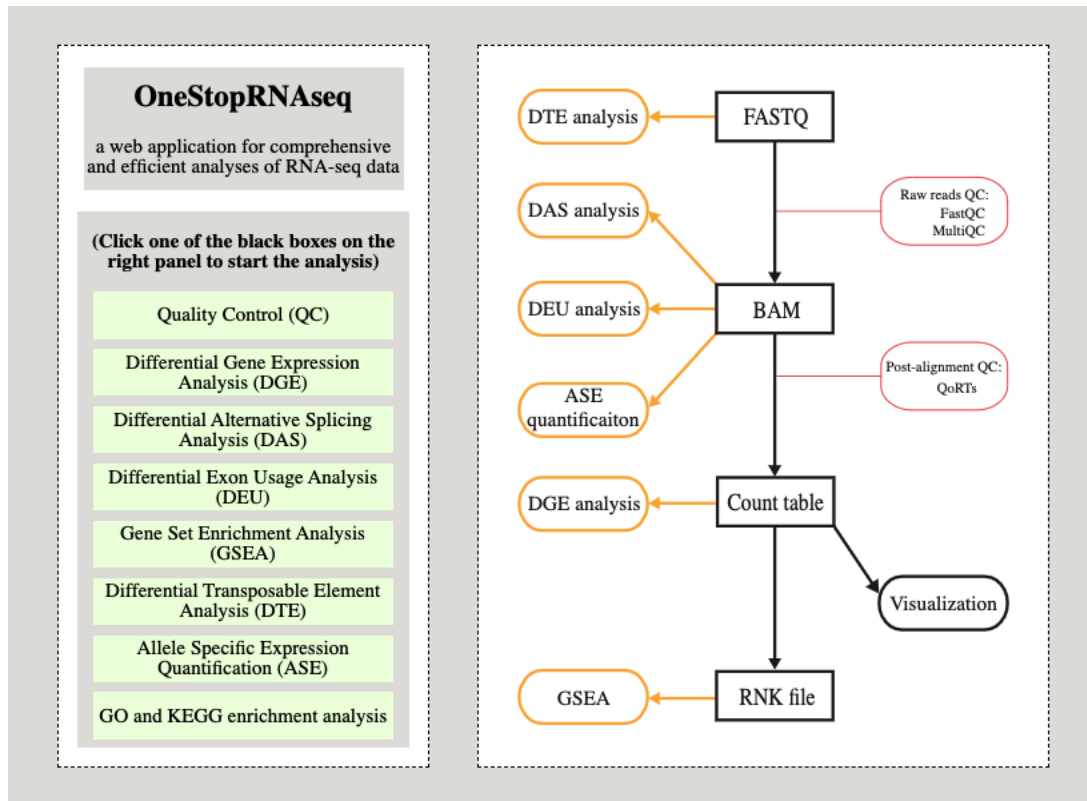


Figure 1. A screenshot showing the homepage of the OneStopRNAseq web application

If users choose to start the analysis with “FASTQ” files, all types of analyses will be performed, though ASE quantification requires users to provide a variant call format (VCF) file containing genotype information. If users start from “BAM” files, DTE analysis will not be performed. If users start the analysis with “Count table”, DAS will not be performed. If users start the analysis with “RNK file”, only GSEA will be performed.

In terms of the “Visualization” module, we are currently hosting DEBrowser, shiny-seq and some light yet handy online plotting tools.

Example

Submit new jobs:

In this section, we will demonstrate how to submit new jobs with “FASTQ” as the analysis entry point. As shown in Figure 2, a new study can be created by clicking the “Upload from Dropbox” or “Upload from GEO” button.

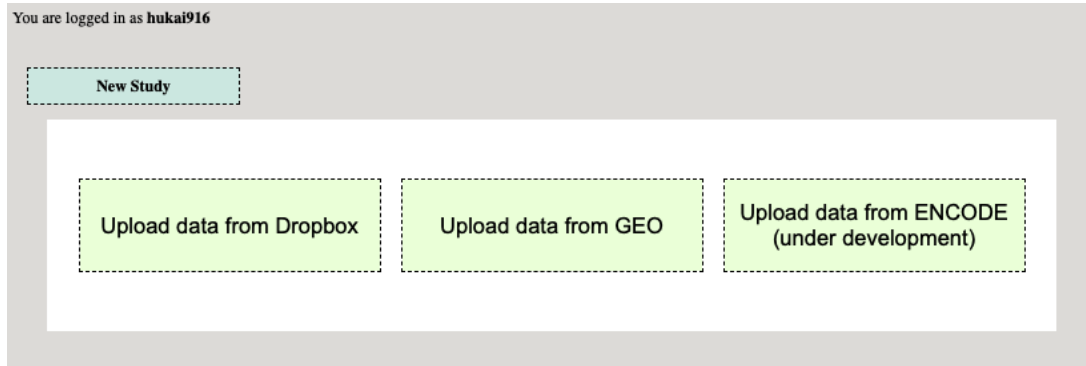


Figure 2. A screenshot showing the web interface for creating a new study

Once “Upload from Dropbox” button is clicked, users will be prompted to specify some basic information about their study. Click “Example 1” to see an example as shown in **Figure 3** below.

Create A New Study

Study name

Reference genome [?]

Sequencing type

Library strandness [?]

Include only uniquely mapped reads [?]

Please specify the number of top gene sets to be plotted [?]

Optionally upload a VCF file (applicable only for ASE):

Figure 3. A screenshot showing the web interface for creating a new study.

Next, users need to provide the metadata (description of the data) and Dropbox links to the data for the study. An example is shown in **Figure 4** below. Alternatively, users can upload an Excel file prefilled with the same information.

New Study: Example_study1

Please enter your data below or upload an Excel file: No file chosen

Group labels should indicate the biological condition of each sample and biological replicates must share the same group label. Otherwise, the analysis results are not valid. Group label examples are control, treated, drug1, drug2, WT, and KO.

By default, batches are set to 1 indicating that all samples are from the same batch. Please make sure batches are assigned correctly and only modify batch numbers if the samples were truly prepared as different batches, e.g., by different researchers, in different dates, or with different batches of reagents. Otherwise, the analysis results are not valid.

Please note that it's not allowed to have any batch assignment that leads to 0 degree of freedom for the error term, e.g., batch 1 assigned to one group and batch 2 to the other group. Every group under comparison must contain all batches.

Group label	Sample label	Batch	R1 file URL	R1 MD5	R2 file URL	R2 MD5	ADD
KO_D0	KO_D0_1_S2	1	https://www.dropbox.com/	7647ab9b2c7bf0d607f	https://www.dropbox.com/	f26255024d96406228b	DELETE
KO_D0	KO_D0_2_S8	1	https://www.dropbox.com/	70d70ce1005e43ef967f	https://www.dropbox.com/	324bfe90f951a60fd785f	DELETE
KO_D2	KO_D2_1_S4	1	https://www.dropbox.com/	d0dd19ba8e6851f5433f	https://www.dropbox.com/	f4b10ba287fdad4b764f	DELETE
KO_D2	KO_D2_2_S10	1	https://www.dropbox.com/	88cddea4d51946a3948f	https://www.dropbox.com/	e994f6eac0abd8554c4c	DELETE
KO_D8	KO_D8_1_S6	1	https://www.dropbox.com/	97bd604ea1236fab89b	https://www.dropbox.com/	ccb994e213c70058189	DELETE
KO_D8	KO_D8_2_S12	1	https://www.dropbox.com/	600c797a9ce7526ff34c	https://www.dropbox.com/	87e5036271269bf8f38f	DELETE
WT_D0	WT_D0_1_S1	1	https://www.dropbox.com/	9d8a903337de135e59ef	https://www.dropbox.com/	44673b1ab77de91902b	DELETE
WT_D0	WT_D0_2_S7	1	https://www.dropbox.com/	ac0f0be56fcb116592a9	https://www.dropbox.com/	4baaf0ac2a5c0ea2ac4c	DELETE
WT_D2	WT_D2_1_S3	1	https://www.dropbox.com/	fab001514ce5016b18cf	https://www.dropbox.com/	72d895b29219dc27bb9	DELETE
WT_D2	WT_D2_2_S9	1	https://www.dropbox.com/	b85d10f59b2425b8c09f	https://www.dropbox.com/	b9ac5611254e20325ed	DELETE
WT_D8	WT_D8_1_S5	1	https://www.dropbox.com/	8acdfaae18ea1f23b678f	https://www.dropbox.com/	1cb8b4a50b814159a8a	DELETE
WT_D8	WT_D8_2_S11	1	https://www.dropbox.com/	1705dc8981e31f680fec	https://www.dropbox.com/	cb1ff48ab4bdb77a0ad0	DELETE

Figure 4. A screenshot showing the web interface to inputting metadata

After clicking the “Submit” button, users can modify the default parameter setting to fit their own needs for DEG and DAS analysis (Figure 5). For example, users can set the minimum log₂ fold change (min_LFC) and maximum false discovery rate (max_FDR) for identifying differentially expressed genes or alternative splicing events. OneStopRNAseq is so flexible that users can make any types of contrasts/comparisons among different sample groups by clicking “Add more” button and selecting treatments to be included in group A and group B. In the results, LFC are computed as $\log_2(\text{groupA} / \text{groupB})$. The following example (Figure 5-1) shows that KO_D0, KO_D2, and KO_D8 are selected for group A, and WT_D0, WT_D2, and WT_D8 are selected for group B, which means you are interested in obtaining the main effect of KO vs WT. If there is a significant interaction between genotype and time points, then different contrasts can be easily specified as following (Figure 5-2). Select KO_D0 for group A and WT_D0 for group B to compare KO and WT at time point 0. Select KO_D2 for group A and WT_D2 for group B to

compare KO and WT at time point 2. Select KO_D8 for group A and WT_D8 for group B to compare KO and WT at time point 8. Users can also set up contrasts easily to compare different time points within each genotype. Select WT_D2 for group A and WT_D0 for group B to compare time 2 and 0 for WT. Select WT_D8 for group A and WT_D0 for group B to compare time 8 and 0 for WT. Similarly, you can set up same contrasts for KO. If you are interested in genes that behaves differently in KO from time point 0 to 2 or 8 comparing to WT, then you simply select KO_D2 or KO_D8 and WT_D0 for group A and WT_D2 or WT_D8 and KO_D0 for group B.

New Study: Example_study1

Please set up the comparisons between GROUPS A and B and the corresponding parameters for the following analyses.
The logFC will be computed as $\log_2(\text{expression of GROUP A} / \text{expression of GROUP B})$.

Differential gene expression analysis using DESeq2

Select sample group(s) for GROUP A: KO_D0 KO_D2 KO_D8 WT_D0 WT_D2 WT_D8

Select sample group(s) for GROUP B: KO_D0 KO_D2 KO_D8 WT_D0 WT_D2 WT_D8

Max_FDR (Maximum False Discovery Rate) [Add more](#)

Min_LFC (Minimum Log2 Fold Change)

Below are optional:

cooksCutoff False True

independentFiltering False True

Alternative splicing and exon usage analysis using rMATS and DEXSeq

Select sample groups for GROUP A: KO_D0 (2 samples) KO_D2 (2 samples) KO_D8 (2 samples) WT_D0 (2 samples) WT_D2 (2 samples) WT_D8 (2 samples)

Select sample group(s) for GROUP B: KO_D0 (2 samples) KO_D2 (2 samples) KO_D8 (2 samples) WT_D0 (2 samples) WT_D2 (2 samples) WT_D8 (2 samples) [Add more](#)

[Go back](#) [Submit](#)

Figure 5-1. A screenshot showing the interface to specify analysis specific parameters and contrasts

New Study: Example_study1

Please set up the comparisons between GROUPS A and B and the corresponding parameters for the following analyses.
The logFC will be computed as $\log_2(\text{expression of GROUP A} / \text{expression of GROUP B})$.

Differential gene expression analysis using DESeq2

Select sample group(s) for GROUP A: <input checked="" type="checkbox"/> KO_D0 <input type="checkbox"/> KO_D2 <input type="checkbox"/> KO_D8 <input type="checkbox"/> WT_D0 <input type="checkbox"/> WT_D2 <input type="checkbox"/> WT_D8	Select sample group(s) for GROUP B: <input type="checkbox"/> KO_D0 <input type="checkbox"/> KO_D2 <input type="checkbox"/> KO_D8 <input checked="" type="checkbox"/> WT_D0 <input type="checkbox"/> WT_D2 <input type="checkbox"/> WT_D8	Max_FDR (Maximum False Discovery Rate) <input type="text" value="0.05"/> Min_LFC (Minimum Log2 Fold Change) <input type="text" value="0.585"/> Below are optional: cooksCutoff ? <input type="radio"/> False <input checked="" type="radio"/> True independentFiltering ? <input checked="" type="radio"/> False <input type="radio"/> True	<input type="button" value="Add more"/>
Select sample group(s) for GROUP A: <input type="checkbox"/> KO_D0 <input checked="" type="checkbox"/> KO_D2 <input type="checkbox"/> KO_D8 <input type="checkbox"/> WT_D0 <input type="checkbox"/> WT_D2 <input type="checkbox"/> WT_D8	Select sample group(s) for GROUP B: <input type="checkbox"/> KO_D0 <input type="checkbox"/> KO_D2 <input type="checkbox"/> KO_D8 <input type="checkbox"/> WT_D0 <input checked="" type="checkbox"/> WT_D2 <input type="checkbox"/> WT_D8		<input type="button" value="Delete"/>
Select sample group(s) for GROUP A: <input type="checkbox"/> KO_D0 <input type="checkbox"/> KO_D2 <input checked="" type="checkbox"/> KO_D8 <input type="checkbox"/> WT_D0 <input type="checkbox"/> WT_D2 <input type="checkbox"/> WT_D8	Select sample group(s) for GROUP B: <input type="checkbox"/> KO_D0 <input type="checkbox"/> KO_D2 <input type="checkbox"/> KO_D8 <input type="checkbox"/> WT_D0 <input type="checkbox"/> WT_D2 <input checked="" type="checkbox"/> WT_D8		<input type="button" value="Delete"/>

Figure 5-2. A screenshot showing interface to set up multiple contrasts to compare genotype effects at different time points

After clicking the “Submit” button, users will be presented a review and verification page before finalizing their submission. After the final submission, your job will be sent to our server. You will receive an email once the job is done. You can also view the analysis status by clicking the “View History” link on the right side of the top menu as detailed below.

View history:

Once you click on the “View History” link, all of your submitted jobs will be listed as shown in **Figure 6** below. Results can be viewed and downloaded if the “Analysis status code” starts with the number 5. A list of available status codes and their interpretation are available on the top of the analysis history. Briefly, code starting with 4 means that job is running, code starting

with 5 indicates that job is finished, and the last 3 digits represent the percent of jobs successfully run. For example, code 5098 means 98% of jobs has been finished, and code 5000 means that the job has been successfully submitted with 0% finished.

You are logged in as hukai916

[View History](#)

Here are a list of possible analysis status codes.

- 0: User's input was received successfully.
- 1: Data was uploaded successfully.
- 2: The configuration was created successfully.
- 3: The Snakemake file is ready to run.
- 4xxx: The job is running, xxx represents the percentage of tasks that have been completed.
- 5xxx: The job is done! xxx indicates the percentage of tasks that have been successfully completed.
- 9xx: The job failed.
 - 900: User's input was not received successfully.
 - 901: Data was not uploaded successfully.
 - 902: The configuration was not created successfully.
 - 903: The Snakemake file is not ready to run.

Study_10: Example_study1_07.59.13-05.07.2020

Analysis name	Analysis status	Analysis result
Comprehensive RNAseq pipeline	4075	Please click here to view the analysis results.

Study_11: Example_study1_08.48.06-05.07.2020

Analysis name	Analysis status	Analysis result
Comprehensive RNAseq pipeline	5097	Please click here to view the analysis results.

Study_12: Example_study1_09.19.50-05.07.2020

Analysis name	Analysis status	Analysis result
Comprehensive RNAseq pipeline	5100	Please click here to view the analysis results.

Figure 6. A screenshot showing the analysis history page

To further choose which type of analysis to view, please click the link below each analysis as shown in **Figure 7**.

Analysis results for the study: Example_study1_08.50.47-08.18.2020

QC metrics
[Raw reads QC by FastQC](#)
[Raw reads QC by MultiQC](#)
[Post-alignment QC by QoRTs](#)
[Download all FastQC files](#) and [Download all MultiQC files](#)

Alignment Results
[CPM normalized bigWig files](#)
[Sorted BAM files](#)

Differential Gene Expression Analysis (DGEA) Results from DESeq2
[Summary](#)
[Download all output files including differentially expressed genes](#)

Gene Set Enrichment Analysis Results from GSEA
[Browse the results](#)
[Download all results](#)

Differential Exon Usage Analysis Results from DEXSeq
[Download all output files including differentially expressed exons](#)

Alternative Splicing Analysis Results from rMATS
[Differential alternative splicing events](#)
[Download all results](#)

Differential Transposon Element Expression Analysis Results from SalmonTE
Results are included in DGEA results. Please note that this analysis is only available for human and mouse with input FASTQ or BAM files.

Figure 7. A screenshot showing the interface for reviewing analysis results

Starting with GEO accession number

To analyze data from GEO, users only need to provide a GEO accession number without the need to download the data to their local computer or dropbox first (**Figure 8**).

Create A New Study

Study name

Reference genome [?]

Sequencing type

Library strandness [?]

Include only uniquely mapped reads [?]

Enter GEO accession number(s) separated by comma [?]

Please specify the number of top gene sets to be plotted [?]

Optionally upload a VCF file (applicable only for ASE):

Figure 8. A screenshot showing the web interface for creating a new study with a GEO accession number.

After users click on the submit button, all associated sample information for the given GEO accession number will be automatically retrieved as shown in **Figure 9**. If multiple GEO accession numbers are provided, please separate them by comma. Please note that the “Group labels” are parsed out from “Sample labels” and manual corrections might be required. The “Batch” values are all set to 1 by default, users

can modify them according to the experimental design. The subsequent steps are the same as starting from Dropbox links.

New Study: Example_study1

Group labels are parsed from the corresponding sample labels automatically to indicate the biological condition of each sample. Biological replicates must share the same group label. Otherwise, the analysis results are not valid. Group label examples are control, treated, drug1, drug2, WT, and KO.

Please make sure batches are assigned correctly and only modify batch numbers if the samples were truly prepared as different batches, e.g., by different researchers, in different dates, or with different batches of reagents. Otherwise, the analysis results are not valid.

Please note that it's not allowed to have any batch assignment that leads to 0 degree of freedom for the error term, e.g., batch 1 assigned to one group and batch 2 to the other group. Every group under comparison must contain all batches.

Group label	Sample label	Batch	GSM ID	SRX ID	ADD
GSE155550_C	GSE155550_C1	1	GSM4706021	SRX8867460	DELETE
GSE155550_C	GSE155550_C2	1	GSM4706022	SRX8867461	DELETE
GSE155550_C	GSE155550_C3	1	GSM4706023	SRX8867462	DELETE
GSE155550_M	GSE155550_M1	1	GSM4706024	SRX8867463	DELETE
GSE155550_M	GSE155550_M2	1	GSM4706025	SRX8867464	DELETE
GSE155550_M	GSE155550_M3	1	GSM4706026	SRX8867465	DELETE

Go back **Submit**

Figure 9. A screenshot showing the metadata retrieved for a given GEO accession number.

Example – other analysis entry points:

Instead of entering from “FASTQ” or “BAM” file, you can also choose to enter the pipeline from “Count table” or “RNK file”. Once you click “Count table” or “RNK file” on the main workflow, you will be prompted to provide required data such as “Count table” as shown in **Figure 10**. Click on the “Example file” button to download the example input files. To start analysis from “Count Table”, see below for detailed file format requirements.

The screenshot shows a web form with the following sections:

- Email address (not required):** A text input field.
- Choose reference genome:** A dropdown menu with the selected option "Mus musculus (select this if using exmaple count file)".
- Upload your count file:** A file upload area with a "Choose File" button, the text "No file chosen", and an "Example file" button.
- Upload your metadata file:** A file upload area with a "Choose File" button, the text "No file chosen", and an "Example file" button.
- Upload your contrast file:** A file upload area with a "Choose File" button, the text "No file chosen", and an "Example file" button.
- Submit:** A green button at the bottom of the form.

Figure 10. A screenshot showing the interface for starting analysis with a count file.

Count file: It is an Excel table with the first column as gene-id/gene-name, column two to six as additional annotation, and other columns as expression quantification as number of reads. Please note that only integers representing the raw read numbers can be used here. If you input normalized expression values, the statistics reported by DESeq2 will be invalid. If you don't have additional annotation, column 2-6 can be left empty.

Metadata file: If you start the workflow from FASTQ or BAM files, metadata such as SAMPLE_LABEL, GROUP_LABEL, and BATCH information can be entered during the data submission. When you start from Count Table, you can upload an Excel file with the metadata filled in using the provided template file meta.xlsx. Please make sure that you label the samples consistently, as labels are case sensitive and different labels are considered different samples, groups, or batches by the analysis workflow. When you start from FASTQ or BAM files, 'SAMPLE_LABEL' should be the same as FASTQ or BAM file prefixes, as OneStopRNAseq with use this to match metadata with input files. For example, 'SAMPLE_LABEL' of KD_D0

should be matched with KD_D0.R1.fastq.gz and KD_D0.R2.fastq.gz for paired-end (PE) RNAseq reads, with KD_D0.fastq.gz for single-end (SE) RNAseq reads, or with KD_D0.bam for aligned BAM file.

Contrast/Comparison file for DGE analysis: It is an Excel file with three rows. Briefly, row 1 contains the header, row 2 contains the treatment groups, row 3 contains the control groups. The resulting log₂ fold change (LFC) are calculated as $\log_2(\text{treatment}/\text{control})$. The group labels should be one of the group labels from meta.xlsx. If you need to include more than one group in the treatment or control group, just separate group labels with semicolon, as shown in the template. You can download the template Excel table contrast.de.xlsx and modify it accordingly.

Contrast/Comparison file for DAS analysis: This file uses the same format and specification as contrast.de.xlsx with one additional constraint, i.e., treatment group and control group need to have the same number of samples per the requirement of rMATS.

To start the analysis from “RNK file”, you will need to specify the “Rank file” as below.

Rank file: It is a tab-delimited file with two columns sorted by the second column in ascending order. The first column should contain the gene-symbol. The second column should be numeric that will be used to rank genes, such as LFC (log fold change), gene expression level, or even binding affinity from ChIP-seq or other experiments. Currently only human and mouse gene-symbols are supported. A template can be downloaded as ‘name1.rnk.txt’.

It is important not to include spaces in the name of any input files or sample/group labels. For sample/group labels, don’t start the label with numbers, and only alphabets, underscore, numbers are accepted.