

# DEGoldS v1.0

Differential Expression analysis pipelines  
benchmarking workflow based on  
Gold-Standard construction

<https://github.com/GGFHF/DEGoldS>

## Table of contents

<b>1. DISCLAIMER</b> .....	3
<b>2. DEPENDENCIES</b> .....	4
Bash scripts (.sh files) .....	4
Miniconda3 .....	4
FastQC (Bioconda installation) .....	4
Trinity (Bioconda installation) .....	5
BUSCO (Bioconda installation) .....	5
QUAST (Bioconda installation) .....	5
GMAP (Bioconda installation) .....	5
Bowtie2 (Bioconda installation) .....	5
NGShelper .....	5
RSEM .....	5
GffCompare .....	5
R scripts (.R files) .....	6
DESeq2 .....	6
edgeR .....	6
Tximport .....	6
readr .....	6
IHW .....	6
Apeglm .....	6
<b>3. WORKFLOW PROCEDURE</b> .....	7
General description .....	7
Gold-standard construction .....	7
Scripts considerations .....	8
<b>4. SCRIPTS DESCRIPTION</b> .....	9
<b>5. VARIABLES INSTRUCTIONS</b> .....	12
In Bash scripts .....	12
In R scripts .....	14

## 1. DISCLAIMER

DEGoldS is available for free download from the GitHub software repository <https://github.com/GGFHF/DEGoldS> under GNU General Public License v3.0. This software has been developed in memoriam of Dr. Pablo G. Goicoechea who initiated this research along with researchers from NEIKER, Departamento de Sistemas y Recursos Naturales (Universidad Politécnica de Madrid) and Universidad del País Vasco (UPV/EHU).

## 2. DEPENDENCIES

### Bash scripts (.sh files)

The following software is necessary to run Bash scripts (.sh files). The scripts are designed to use the mentioned programs from the Bioconda repositories (using for example Miniconda3); however, any other installation should work as long as the software (and its dependencies) directories are included in the PATH variable. In this case the following lines in the scripts could be commented (#) to avoid errors and warnings:

```
#cd "$CONDA_DIR"
```

```
#source activate "$FASTQC_ENV"
```

And the following line should be added to include the directory hosting the desired program in the PATH variable:

```
PATH=/PATH_TO_YOUR_APP_DIRECTORY:$PATH
```

The instructions to install the software through Bioconda repositories using Miniconda are as follows

### Miniconda3

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY (where you will install Miniconda3)
```

```
$ wget https://repo.continuum.io/miniconda/Miniconda3-latest-Linux-x86_64.sh
```

```
$ chmod u+x Miniconda3-latest-Linux-x86_64.sh
```

```
$ ./Miniconda3-latest-Linux-x86_64.sh -b -p /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3
```

```
$ rm Miniconda3-latest-Linux-x86_64.sh
```

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
```

```
$ ./conda config --add channels defaults
```

```
$ ./conda config --add channels conda-forge
```

```
$ ./conda config --add channels r
```

```
$ ./conda config --add channels bioconda
```

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
```

```
$ ./conda install --yes --name base mamba
```

### FastQC (Bioconda installation)

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
```

```
$ ./mamba create --yes --name fastqc fastqc
```

**Trinity (Bioconda installation)**

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
$ ./mamba create --yes --name trinity trinity
```

**BUSCO (Bioconda installation)**

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
$ ./mamba create --yes --name busco busco
```

**QUAST (Bioconda installation)**

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
$ ./mamba create --yes --name quast quast
```

**GMAP (Bioconda installation)**

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
$ ./mamba create --yes --name gmap gmap
```

**Bowtie2 (Bioconda installation)**

```
$ cd /PATH_TO_YOUR_APPS_DIRECTORY/Miniconda3/bin
$ ./mamba create --yes --name bowtie2 bowtie2
```

**NGShelper**

Installation instruction in <https://github.com/GGFHF/NGShelper>. These instructions don't make NGShelper available from PATH variable. However, this is not necessary as software directory has to be specified in NGSHELPER\_DIR variable in the scripts.

**RSEM**

Although RSEM can be installed also from Bioconda repositories, issues regarding to reads simulations have been experienced. For this reason, installation from source code is suggested (instruction in <https://github.com/deweylab/RSEM>). These instructions make RSEM available from PATH variable, necessary to run the scripts as they are.

**GffCompare**

Although GffCompare can be installed also from Bioconda repositories some issues have been experienced when running. For this reason, installation from source code is suggested (instruction in <https://github.com/gpertea/gffcompare>). These instructions don't make GffCompare available from PATH variable. However, this is not necessary as software directory has to be specified in GFFCOMPARE\_DIR variable in the scripts.

### R scripts (.R files)

The following R packages are necessary to run R scripts (.R files). All scripts have been tested in R version 4.

#### DESeq2

```
if (!require("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("DESeq2")
```

#### edgeR

```
if (!require("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("edgeR")
```

#### Tximport

```
if (!require("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("tximport")
```

#### readr

```
install.packages("readr")
```

#### IHW

```
if (!require("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("IHW")
```

#### Apeglm

```
if (!require("BiocManager", quietly = TRUE))  
  install.packages("BiocManager")  
BiocManager::install("apeglm")
```

### 3. WORKFLOW PROCEDURE

#### General description

DEGoldS (Differential Expression analysis pipelines benchmarking workflow based on Gold-Standard construction) is divided into 4 main steps which are divided at the same time into different substeps consisting on different scripts. It is suggested to create a different directory for every script during the workflow. This simulation-based workflow is designed for testing DE pipelines using two group of samples and 3 samples per group starting from real reads so results can be optimized for a particular set of samples. Although any DE analyser could be used for benchmarking this workflow pretends to use DESeq2 and edgeR, so in case user wants to test any other new scripts should be used. This workflow requires that the pipeline being tested outputs a GTF/GFF assembly (e. g. Stringtie) so transcripts IDs conversion can be done. If this is not possible a transcriptome to genome alignment with GMAP is suggested (e. g. script 1d-03) as it can output results in GFF file.

Two different simulation pipelines covering from a simpler and less realistic simulation (Sim1) to a more realistic but more complex one (Sim2) are described. Both procedures share all the steps except Step 2b as Sim2 gets expression values from real data while Sim1 has unrealistic expression values with no dispersion between samples. Step 2b consists on 4 modules: A-selectable transcripts pool construction, B-helper scripts for selectable counts pool matrix construction, C-Sim1 TPM tables construction and D-Sim2 TPM tables construction. "A" module is used before either "C" (Sim1) or "D" (Sim2) modules, however "B" module is only used before "D" module. "B" module together with 2b-D01, 2b-D02, 2b-D03 and 2b-D04 helps to get a reliable gold-standard.

#### Gold-standard construction

The gold-standard is built from differential gene expression analysis (scripts 2d) of the count values taken from script 2c-01. At this point the gold-standard is called "simulation gold-standard" (simGS) because gene and transcripts IDs belong to the transcriptome where simulated reads come from. On the contrary, when this IDs are converted into IDs from the assembly of the tested DE pipeline (script 4a-02) is called "pipeline gold-standard" (pipelineGS). A reliable or suitable gold-standard for benchmarking means that simGS contains most of the "upregulated" transcripts requested in simulation but does not contain any "filler" transcripts as these should not show any significant differences in a DE expression analysis. Using real expression values (i.e. Sim2) could lead into a gold-standard with upregulated genes other than the desired ones or hiding those desired ones. Provided scripts optimize the chances (so the times to repeat the procedure) to get a count matrix with potential interesting real counts for upregulated and filler transcripts for DESeq2 and edgeR. However, after every simulation simGS should be checked to make sure that most of the genes corresponding to upregulated transcripts are recovered while no genes corresponding to filler transcripts are in it (scripts 2d). The rationale of this step consists on removing counts from a real matrix that could be in the limit between significant and non-significant differentially expressed values so increasing chances for getting a suitable gold-standard. The only thing the user has to do is to define the number of sequences (line of count in a matrix) to remove up and down significance limit according to the total number of transcripts within the transcriptome (script 2b-D02). This would reduce the times a simulation has to be performed until a suitable gold-standard is reached.

It is important to make sure that when converting simGS into pipelineGS there are no multiple pipeline genes corresponding to simulated “upregulated” transcripts in file “03-DESeq2\_simGS\_correspondence\_table.csv” in script 4a-02. This should not be usual but in case this happens user should decide how to proceed with benchmarking process. One solution could be to consider all the possible genes but counting +1TP (True Positive) in the scoring process.

### Scripts considerations

The scripts have to be only modified in the section "Variables definition" with the required strings provided. The strings should be placed between quotes ("""). Usually strings are characters strings, but if specified in instructions a number should be placed. If there are no quotes it means a number should be placed. Explanation of every variable can be seen in section 5 (variables instructions) as well as in the upper part of every script.



## 4. SCRIPTS DESCRIPTION

### **Step 1 -> Transcriptome assembly and alignment to reference genome**

Script 1a-01: performs QC on input reads with FASTQC

Script 1b-01: assembles input reads into a de novo transcriptome with Trinity

Script 1b-02: filtrates de novo transcripts by lenght with NGShelper

Script 1c-01: assess transcriptome quality with QUAST

Script 1c-02: assess transcriptome completeness with BUSCO

Script 1d-01: builds reference genome index for GMAP

Script 1d-02: aligns transcriptome against reference genome with GMAP

Script 1d-03: aligns transcriptome against reference genome with GMAP (GFF output)

Script 1d-04: classifies de novo transcriptome transcripts into 0 times alignes (path 0), uniquely aligned (path 1) and multiple times aligned (path n) transcripts

Script 1d-05: modifies GFF output from GMAP for unstranded transcriptomes

### **Step 2 -> Simulation of reads files**

Script 2a-01: performs learning step from input reads and de novo transcriptome with RSEM

Script 2b-A01: classifies transcripts according to locus in the reference genome

Script 2b-A02: gets transcripts pools for "upregulated" and "filler" transcripts selection

Script 2b-B01: transcripts DE analysis with DESeq2 from learnrign step data

Script 2b-B02: transcripts DE analysis with DESeq2 from learnrign step data with no pre-filtering step

Script 2b-B03: transcripts DE analysis with edgeR from learnrign step data

Script 2b-B04: transcripts DE analysis with edgeR from learnrign step data with no pre-filtering step

Script 2b-B05: classifies common/not common significant transcripts in DE analysis in scripts 2b-B01 and 2b-B02.

Script 2b-C01: selects upregulated transcripts (only for Sim1)

Script 2b-C02: checks that upregulated transcripts belong to different genes (only for Sim1)

Script 2b-C03: selects filler transcripts and builds the new TPM tables for reads simulation (only for Sim1)

Script 2b-C04: performs reads simulation with RSEM (only for Sim1)

Script 2b-D01: gets count matrix building from 2a-01 output (only for Sim2)

Script 2b-D02: builds the selectable count pool matrix (only for Sim2)

Script 2b-D03: selects significant and not significant counts (only for Sim2)

Script 2b-D04: builds the significant and not significant count matrix (only for Sim2)

Script 2b-D05: selects upregulated transcripts (only for Sim2)

Script 2b-D06: checks that upregulated transcripts belong to different genes (only for Sim2)

Script 2b-D07: selects filler transcripts (only for Sim2)

Script 2b-D08: checks that filler transcripts belong to different genes (only for Sim2)

Script 2b-D09: obtains transcripts from transcriptome that will have 0 counts -surplus transcripts- (only for Sim2)

Script 2b-D10: prepare counts tables (isoforms.results files) for 2b-D11 (only for Sim2)

Script 2b-D11: builds the new TPM tables for reads simulation (only for Sim2)

Script 2b-D12: performs reads simulation with RSEM (only for Sim2)

Script 2c-01: generates simulated reads from de novo transcriptome with RSEM

Script 2d-01: performs DE analysis from simulation stats to build simulation gold-standard with DESeq2 (simGS)

Script 2d-02: performs DE analysis from simulation stats to build simulation gold-standard with edgeR (simGS)

### **Step 3 -> To be tested pipeline running**

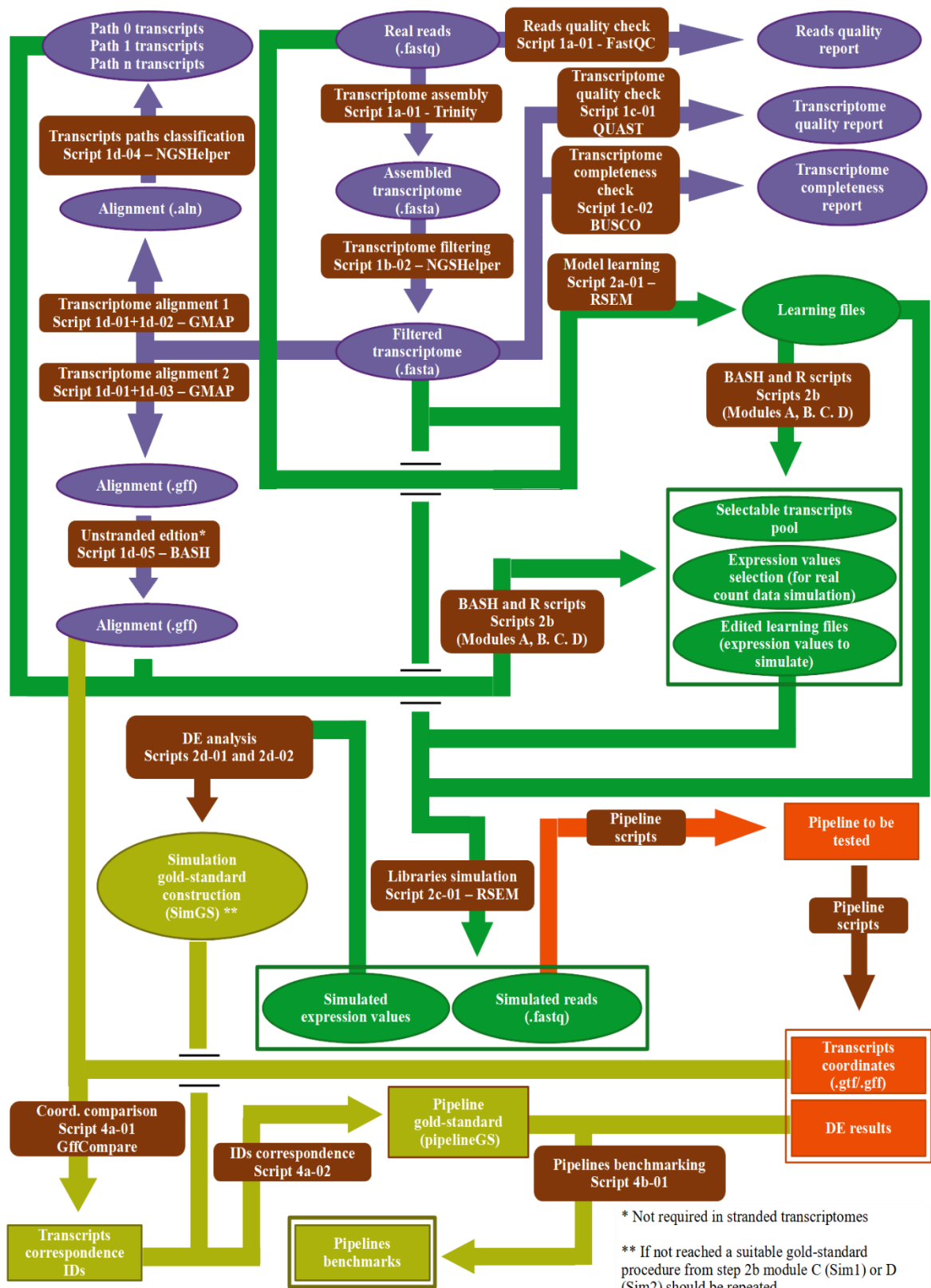
In this step each user should include its own scripts from the DE pipeline to compare.

### **Step 4 -> Benchmarking of the tested DE analysis pipeline**

Script 4a-01: gets coordinates correspondence between de novo transcripts and transcripts assembled in tested pipeline using GFF output from GMAP alignment and GTF assembly annotation using GffCompare

Script 4a-02: outputs pipeline gold-standard (pipelineGS) from the correspondance in script 4a-01 and simGS in Step 2d

Script 4b-01: performs benchmarking of the pipeline according to pipelineGS and results in DE analysis in the tested pipeline



## 5. VARIABLES INSTRUCTIONS

### In Bash scripts

BOWTIE2\_DIR: Bowtie2 directory

BUSCO\_ENV: BUSCO environment in Conda

COMMONDEIDS: "03-DESeq2\_edgeR\_common\_sig\_isoforms.csv" file path from script 2b-B05

CONDA\_DIR: Conda directory

CORRESP\_FILE\_PATH: "ref\_pipeline.txt" file path from 4a-01

COUNTMATRIX: "01-learn\_count\_matrix.csv" file path from script 2b-D01

DE\_RES\_FILE\_DEL: column delimiter in \$DE\_RES\_FILE\_PATH

DE\_RES\_FILE\_GENE\_ID\_COL\_NUMBER: gene ID column number in \$DE\_RES\_FILE\_PATH

DE\_RES\_FILE\_PATH: DE analysis results table from tested pipeline

DESEQ2\_LEARN\_DE\_ALL\_RES: "01-learning\_DESeq2\_all\_seqs.csv" file path from script 2b-B01

DESEQ2\_LEARN\_DE\_SIG\_RES: "02-learning\_DESeq2\_sig\_seqs.csv" file path from script 2b-B01

EDGER\_LEARN\_DE\_ALL\_RES: "01-learning\_edgeR\_all\_seqs.csv" file path from script 2b-B01

EDGER\_LEARN\_DE\_SIG\_RES: "02-learning\_edgeR\_sig\_seqs.csv" file path from script 2b-B03

EDITED\_GMAP\_GFF\_FILE: edited GMAP GFF output (unstranded\_edited\_alignment.gff in 1d-05)

EDITED\_SAMPLES\_LEARN\_DIR: directory with generated "isoforms.results" in script 2b-C03 (Sim1) or 2b-D11 (Sim2)

FASTQC\_ENV: FASTQC environment in Conda

FILLER\_SEL\_SEQS\_FILE\_PATH: "filler\_seqs.txt" file path from script 2b-D07

FILTNUMNOSIGDESEQ2: number of DESeq2 non-significant sequences to be filtered starting from head

FILTNUMNOSIGEDGER: number of edgeR non-significant sequences to be filtered starting from head

FILTNUMSIGDESEQ2: number of DESeq2 significant sequences to be filtered starting from bottom

FILTNUMSIGEDGER: number of edgeR significant sequences to be filtered starting from bottom

GENOME\_FILE\_PATH: gzipped reference genome FASTA file path

GFFCOMPARE\_DIR: GffCompare directory

GMAP\_ALIGNMENT\_FILE: (alignment.aln in 1d-02)

GMAP\_ENV: GMAP environment in Conda

GMAP\_GFF\_FILE: GMAP GFF output (alignment.gff in 1d-03)

GMAP\_INDEX\_DIR: reference genome GMAP index directory  
 GMAP\_INDEX\_NAME: reference genome GMAP index name  
 GROUP01\_UPREG\_SEQS\_X1\_FILE\_PATH: "group01\_upreg\_seqs\_x1.txt" file path from 2b-C01  
 GROUP01\_UPREG\_SEQS\_X2\_FILE\_PATH: "group01\_upreg\_seqs\_x2.txt" file path from 2b-C01  
 GROUP01\_UPREG\_SEQS\_X3\_FILE\_PATH: "group01\_upreg\_seqs\_x3.txt" file path from 2b-C01  
 GROUP01\_UPREG\_SEQS\_X4\_FILE\_PATH: "group01\_upreg\_seqs\_x5.txt" file path from 2b-C01  
 GROUP01\_UPREG\_SEQS\_X5\_FILE\_PATH: "group01\_upreg\_seqs\_x4.txt" file path from 2b-C01  
 GROUP02\_UPREG\_SEQS\_X1\_FILE\_PATH: "group01\_upreg\_seqs\_x1.txt" file path from 2b-C01  
 GROUP02\_UPREG\_SEQS\_X2\_FILE\_PATH: "group02\_upreg\_seqs\_x2.txt" file path from 2b-C01  
 GROUP02\_UPREG\_SEQS\_X3\_FILE\_PATH: "group03\_upreg\_seqs\_x3.txt" file path from 2b-C01  
 GROUP02\_UPREG\_SEQS\_X4\_FILE\_PATH: "group04\_upreg\_seqs\_x4.txt" file path from 2b-C01  
 GROUP02\_UPREG\_SEQS\_X5\_FILE\_PATH: "group05\_upreg\_seqs\_x5.txt" file path from 2b-C01  
 GS\_NAME: gold-standard given name (recommended a name distinguishing each DE analysers, i.e. DESeq2/edgeR)  
 LEARN\_DIR: "02-model\_learning/output/" directory from script 2a-01  
 LEFT\_SUFFIX: left reads file suffix (including extension)  
 LINEAGE\_DATASET: specify the name of the BUSCO lineage to be used (BUSCO MANUAL --lineage\_dataset option)  
 LOCI\_FILE\_PATH: "transcriptome.loci" file from script 2b-A01  
 MEMORY: maximum number of memory to use (in GB)  
 NCPU: number of CPUs  
 NF\_DESEQ2\_LEARN\_DE\_SIG\_RES: "02-nf\_learning\_DESeq2\_sig\_seqs.csv" file path from script 2b-B01  
 NF\_EDGER\_LEARN\_DE\_SIG\_RES: "02-learning\_edgeR\_sig\_seqs.csv" file path from script 2b-B01  
 NGSHELPER\_DIR: NGSHelper directory  
 NO\_SIG\_COUNT\_MATRIX\_FILE\_PATH: "06-no\_sig\_rest\_removed\_count\_matrix.csv" file path from script 2b-D02  
 NOTCOMMONDEIDS: "06-DESeq2\_edgeR\_not\_common\_sig\_isoforms.csv" file path from script 2b-B05  
 PATH1\_SEQS\_IDS\_FILE\_PATH: "assembly-ids-1path.txt" from script 1d-04  
 PIPELINE\_ANN\_FILE: tested pipeline GTF/GFF file  
 PIPELINE\_TX2GENE\_FILE\_PATH: two column text file (transcript ID-gene ID) separated by space (view tx2gene.csv example file) for pipeline gene IDs

PIPELINEGS\_FILE\_PATH: "05-\$GS\_NAME\_pipelineGS\_geneIDs.txt" file path from script 4a-02

QUAST\_ENV: QUAST environment in Conda

RAW\_TRANSCRIPTOME\_FILE\_PATH: pre-filtered transcriptome FASTA file path

READS\_DIR: Reads directory

READSNUMBER: number of reads to be simulated

REPNUMBER: number of replicates to be simulated for each input "isoforms.results" file

RIGHT\_SUFFIX: right reads file suffix (including extension)

SAMPLE\_LEARN\_FILE\_PATH: one of the ".isoform.result" files file path in 2a-01

SEL\_NO\_SIG\_SEQS\_FILE\_PATH: "02-sel\_no\_sig\_seqs.txt" file path from script 2b-D03

SEL\_SIG\_SEQS\_FILE\_PATH: "01-sel\_sig\_seqs.txt" file path from script 2b-D03

SIG\_COUNT\_MATRIX\_FILE\_PATH: "05-common\_sig\_rest\_removed\_count\_matrix.csv" file path from script 2b-D02

SIMGS\_FILE\_DEL: column delimiter in \$SIMGS\_FILE\_PATH

SIMGS\_FILE\_GENE\_ID\_COL\_NUMBER: gene ID column number in \$SIMGS\_FILE\_PATH

SIMGS\_FILE\_PATH: "simGS\_DESeq2.csv" (DESeq2) and "simGS\_edgeR.csv" file path (edgeR) from 2d-01 and 2d-02 respectively or any other used

SURPLUS\_OBT\_SEQS\_FILE\_PATH: "surplus\_seqs\_ids.txt" file path from script 2b-D09

TRANSCRIPTOME\_FILE\_PATH: by length filtered transcriptome FASTA file path

TRANSCRIPTOME\_INDEX\_DIR: transcriptome RSEM index directory

TRANSCRIPTOME\_INDEX\_NAME: transcriptome RSEM index name

TRANSCRIPTOME\_NAME: filtered transcriptome name (file name without extension)

TRINITY\_ENV: Trinity environment in Conda

UPREG\_SEL\_SEQS\_FILE\_PATH: "01-DE\_trinity\_transcripts.txt" file path from 2b-C02 (Sim1) and "upreg\_seqs.txt" file path from script 2b-D05 (Sim2)

WD: Working directory

### In R scripts

base\_value: base expression level value

candidate\_upreg\_seqs\_file\_path: a text file with a list of selectable transcript IDs in the first column with no header for upregulated transcripts selection (i. e. selectable transcripts pool), for example, "08-path1\_candidate\_unique\_seqs\_IDs.txt" (Sim1) and "11-path1\_candidate\_unique-plus\_seqs.txt" (Sim2) file path from 2b-A02

counts\_dir: a directory where .sim.isoforms.results from script 2c-01 are placed after changing 5th column name into "expected\_count" instead of "count"

filler\_group01\_sample01\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 1 in group 1

filler\_group01\_sample02\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 2 in group 1

filler\_group01\_sample03\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 3 in group 1

filler\_group02\_sample01\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 1 in group 2

filler\_group02\_sample02\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 2 in group 2

filler\_group02\_sample03\_learn\_file\_path: filler ".isoform.result" file path in 2b-D10 for sample 3 in group 2

filler\_pre\_candidate\_seqs\_file\_path: a text file with a list of selectable transcript IDs in the first column with no header for filler transcripts selection (i. e. selectable transcripts pool), for example, "10-candidate\_unique-plus\_seqs.txt" file path in script 2b-A02

filler\_sel\_seqs\_number: number of filler transcripts to select

group01\_sample01\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 1 in group 1

group01\_sample01\_samplename: sample name for sample 1 in group 1 (should coincide with the names of the input reads in the workflow)

group01\_sample02\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 2 in group 1

group01\_sample02\_samplename: sample name for sample 2 in group 1 (should coincide with the names of the input reads in the workflow)

group01\_sample03\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 3 in group 1

group01\_sample03\_samplename: sample name for sample 3 in group 1 (should coincide with the names of the input reads in the workflow)

group01\_upreg\_seqs\_x1\_number: number of transcripts to select for group 1 with level 1 of expression in upregulated transcripts

group01\_upreg\_seqs\_x2\_number: number of transcripts to select for group 1 with level 2 of expression in upregulated transcripts

group01\_upreg\_seqs\_x3\_number: number of transcripts to select for group 1 with level 3 of expression in upregulated transcripts

group01\_upreg\_seqs\_x4\_number: number of transcripts to select for group 1 with level 4 of expression in upregulated transcripts

group01\_upreg\_seqs\_x5\_number: number of transcripts to select for group 1 with level 5 of expression in upregulated transcripts

group01\_upreg\_X1\_sel\_seqs\_file\_path: selected transcripts in group 1 with level 1 of expression in upregulated transcripts ("group01\_upreg\_seqs\_x1.txt" file path from script 2b-C01)

group01\_upreg\_X2\_sel\_seqs\_file\_path: selected transcripts in group 1 with level 2 of expression in upregulated transcripts ("group01\_upreg\_seqs\_x2.txt" file path from script 2b-C01)

group01\_upreg\_X3\_sel\_seqs\_file\_path: selected transcripts in group 1 with level 3 of expression in upregulated transcripts ("group01\_upreg\_seqs\_x3.txt" file path from script 2b-C01)

group01\_upreg\_X4\_sel\_seqs\_file\_path: selected transcripts in group 1 with level 4 of expression in upregulated transcripts ("group01\_upreg\_seqs\_x4.txt" file path from script 2b-C01)

group01\_upreg\_X5\_sel\_seqs\_file\_path: selected transcripts in group 1 with level 5 of expression in upregulated transcripts ("group01\_upreg\_seqs\_x5.txt" file path from script 2b-C01)

group02\_sample01\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 1 in group 2

group02\_sample01\_samplename: sample name for sample 1 in group 1 (should coincide with the names of the input reads in the workflow)

group02\_sample02\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 2 in group 2

group02\_sample02\_samplename: sample name for sample 2 in group 2 (should coincide with the names of the input reads in the workflow)

group02\_sample03\_learn\_file\_path: ".isoform.result" file in 2a-01 for sample 3 in group 2

group02\_sample03\_samplename: sample name for sample 3 in group 3 (should coincide with the names of the input reads in the workflow)

group02\_upreg\_seqs\_x1\_number: number of transcripts to select for group 2 with level 1 of expression in upregulated transcripts

group02\_upreg\_seqs\_x2\_number: number of transcripts to select for group 2 with level 1 of expression in upregulated transcripts

group02\_upreg\_seqs\_x3\_number: number of transcripts to select for group 2 with level 1 of expression in upregulated transcripts

group02\_upreg\_seqs\_x4\_number: number of transcripts to select for group 2 with level 1 of expression in upregulated transcripts

group02\_upreg\_seqs\_x5\_number: number of transcripts to select for group 2 with level 1 of expression in upregulated transcripts

group02\_upreg\_X1\_sel\_seqs\_file\_path: selected transcripts in group 2 with level 1 of expression in upregulated transcripts ("group02\_upreg\_seqs\_x1.txt" file path from script 2b-C01)

group02\_upreg\_X2\_sel\_seqs\_file\_path: selected transcripts in group 2 with level 2 of expression in upregulated transcripts ("group02\_upreg\_seqs\_x2.txt" file path from script 2b-C01)

group02\_upreg\_X3\_sel\_seqs\_file\_path: selected transcripts in group 2 with level 3 of expression in upregulated transcripts ("group02\_upreg\_seqs\_x3.txt" file path from script 2b-C01)

group02\_upreg\_X4\_sel\_seqs\_file\_path: selected transcripts in group 2 with level 4 of expression in upregulated transcripts ("group02\_upreg\_seqs\_x4.txt" file path from script 2b-C01)

group02\_upreg\_X5\_sel\_seqs\_file\_path: selected transcripts in group 2 with level 5 of expression in upregulated transcripts ("group02\_upreg\_seqs\_x5.txt" file path from script 2b-C01)



learn\_dir: "02-model\_learning/output/" directory from script 2a-01

no\_sig\_count\_matrix\_file\_path: "06-no\_sig\_rest\_removed\_count\_matrix.csv" file path from script 2b-D02

sample\_learn\_file\_path: one of the ".isoform.result" files file path in 2a-01

samples\_description\_file\_path: samples description text file path (view samples.txt example file)

sel\_no\_sig\_count\_matrix\_file\_path: "02-complete\_no\_sig\_count\_matrix.csv" file path from script 2b-D04

sel\_no\_sig\_seqs\_number: number of sequences to select from count matrix in "06-no\_sig\_rest\_removed\_count\_matrix.csv" from script 2b-D02 that will belong to filler transcripts

sel\_sig\_count\_matrix\_file\_path: "01-complete\_sig\_count\_matrix" file path from script 2b-D04

sel\_sig\_seqs\_number: number of sequences to select from count matrix in "05-common\_sig\_rest\_removed\_count\_matrix.csv" from script 2b-D02 that will belong to upregulated transcripts

sig\_count\_matrix\_file\_path: "05-common\_sig\_rest\_removed\_count\_matrix.csv" file path from script 2b-D02

sim\_tx2gene\_file\_path: two column text file (transcript ID-gene ID) separated by space (view tx2gene.csv example file) for transcriptome gene IDs

surplus\_group01\_sample01\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 1 in group 1

surplus\_group01\_sample02\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 2 in group 1

surplus\_group01\_sample03\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 3 in group 1

surplus\_group02\_sample01\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 1 in group 2

surplus\_group02\_sample02\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 2 in group 2

surplus\_group02\_sample03\_learn\_file\_path: surplus ".isoform.result" file path in 2b-D10 for sample 3 in group 2

upreg\_group01\_sample01\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 1 in group 1

upreg\_group01\_sample02\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 2 in group 1

upreg\_group01\_sample03\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 3 in group 1

upreg\_group02\_sample01\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 1 in group 2

upreg\_group02\_sample02\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 2 in group 2

upreg\_group02\_sample03\_learn\_file\_path: upregulated ".isoform.result" file path in 2b-D10 for sample 3 in group 2

upreg\_sel\_seqs\_file\_path: "01-DE\_trinity\_transcripts.txt" file path from script 2b-C02 (Sim1) or "upreg\_seqs.txt" file path from script 2b-D05 (Sim2)

upreg\_seqs\_number: number of upregulated transcripts to select

upreg\_value\_X1: upregulated expression value for level 1

upreg\_value\_X2: upregulated expression value for level 2

upreg\_value\_X3: upregulated expression value for level 3

upreg\_value\_X4: upregulated expression value for level 4

upreg\_value\_X5: upregulated expression value for level 5

wd: Working directory