

Docker image: rMATS-turbo- 0.1

Prerequisite

- Docker

Software installed in this image

- Operating system: Debian GNU/Linux 8 (jessie)
- gcc version 4.9.2 (Debian 4.9.2-10)
- STAR-2.5.2b
- Python 2.7.12
 - Cython (0.25.2)
 - numpy (1.12.0)
- libblas-dev 1.2.20110419-10
- liblapack-dev 3.5.0-4
- libgsl0ldbl 1.16+dfsg-2

install the image

```
1 | docker load -i rmats-turbo-0.1.tar
```

run the image

```
1 | docker run rmats:turbo01 [options]
```

RMATS USAGE

About

rMATS-turbo is the C/Cython version of rMATS (refer to http://rnaseq-mats.sourceforge.net/user_guide.htm). The main difference between rMATS-turbo and rMATS is speed and space usage. The speed of rMATS-turbo is 100 times faster and the output file is 1000 times smaller than rMATS. These advantages make analysis and storage of large scale dataset easy and convenient.

	Counting part	Statistical part
Speed (C/Cython version vs Python version)	20~100 times faster (one thread)	300 times faster (6 threads)
Storage usage (C/Cython version vs Python version)	1000 times smaller	-

Usage

```
1  docker run rmats:turbo01 -h
2
3  usage: usage: rmats.py [options] arg1 arg2
4
5  optional arguments:
6      -h, --help            show this help message and exit
7      --version            Version.
8      --gtf GTF            An annotation of genes and transcripts in GTF format.
9      --b1 B1              BAM configuration file.
10     --b2 B2              BAM configuration file.
11     --s1 S1              FASTQ configuration file.
12     --s2 S2              FASTQ configuration file.
13     --od OD              output folder of post step.
14     -t {paired,single}   readtype, single or paired.
15     --libType {fr-unstranded,fr-firststrand,fr-secondstrand}
16                          Library type. Default is unstranded (fr-unstranded).
17                          Use fr-firststrand or fr-secondstrand for strand-
18                          specific data.
19     --readLength READLENGTH
20                          The length of each read.
21     --anchorLength ANCHORLENGTH
22                          The anchor length. (default is 1.)
23     --tophatAnchor TOPHATANCHOR
24                          The "anchor length" or "overhang length" used in the
25                          aligner. At least "anchor length" NT must be
26                          mapped to each end of a given junction. The default is
27                          1. (This parameter applies only if using fastq).
28     --bi BINDEX           The folder name of the STAR binary indexes (i.e., the
29                          name of the folder that contains SA file). For
30                          example, use ~/STARindex/hg19 for hg19. (Only if using
31                          fastq)
32     --nthread NTHREAD    The number of thread. The optimal number of thread
33                          should be equal to the number of CPU core.
34     --tstat TSTAT        the number of thread for statistical model.
35     --cstat CSTAT        The cutoff splicing difference. The cutoff used in the
36                          null hypothesis test for differential splicing. The
37                          default is 0.0001 for 0.01% difference. Valid: 0 ≤
38                          cutoff < 1.
39     --statoff            Turn statistical analysis off.
```

Output

--od read count generated by the post step:

- fromGTF.AS_Event.txt: all possible alternative splicing (AS) events derived from GTF and RNA.
- JC.raw.input.AS_Event.txt evaluates splicing with only reads that span splicing junctions
 - IJCSAMPLE1: inclusion junction counts for SAMPLE_1, replicates are separated by comma
 - SJCSAMPLE1: skipping junction counts for SAMPLE_1, replicates are separated by comma
 - IJCSAMPLE2: inclusion junction counts for SAMPLE_2, replicates are separated by comma
 - SJCSAMPLE2: skipping junction counts for SAMPLE_2, replicates are separated by comma
 - IncFormLen: length of inclusion form, used for normalization
 - SkipFormLen: length of skipping form, used for normalization
- JCEC.raw.input.AS_Event.txt evaluates splicing with reads that span splicing junctions and reads on target (striped)

regions on home page figure)

- ICSAMPLE1: inclusion counts for SAMPLE_1, replicates are separated by comma
- SCSAMPLE1: skipping counts for SAMPLE_1, replicates are separated by comma
- ICSAMPLE2: inclusion counts for SAMPLE_2, replicates are separated by comma
- SCSAMPLE2: skipping counts for SAMPLE_2, replicates are separated by comma
- IncFormLen: length of inclusion form, used for normalization
- SkipFormLen: length of skipping form, used for normalization
- AS_Event.MATS.JC.txt evaluates splicing with only reads that span splicing junctions
 - ICSAMPLE1: inclusion counts for SAMPLE_1, replicates are separated by comma
 - SCSAMPLE1: skipping counts for SAMPLE_1, replicates are separated by comma
 - ICSAMPLE2: inclusion counts for SAMPLE_2, replicates are separated by comma
 - SCSAMPLE2: skipping counts for SAMPLE_2, replicates are separated by comma
- AS_Event.MATS.JCEC.txt evaluates splicing with reads that span splicing junctions and reads on target (striped regions on home page figure)
 - ICSAMPLE1: inclusion counts for SAMPLE_1, replicates are separated by comma
 - SCSAMPLE1: skipping counts for SAMPLE_1, replicates are separated by comma
 - ICSAMPLE2: inclusion counts for SAMPLE_2, replicates are separated by comma
 - SCSAMPLE2: skipping counts for SAMPLE_2, replicates are separated by comma
- Important columns contained in output files above
 - IncFormLen: length of inclusion form, used for normalization
 - SkipFormLen: length of skipping form, used for normalization
 - P-Value: (The meaning of p value???)
 - FDR: (The meaning of FDR???)
 - IncLevel1: inclusion level for SAMPLE_1 replicates (comma separated) calculated from normalized counts
 - IncLevel2: inclusion level for SAMPLE_2 replicates (comma separated) calculated from normalized counts
 - IncLevelDifference: $\text{average}(\text{IncLevel1}) - \text{average}(\text{IncLevel2})$
- bamX_Y STAR mapping result.

How to transfer data to docker image's file system.

Docker has it's own file system, called Union File System. We're not going to dig into these concepts. What we're going to do is to learn how we can manage data inside and between our Docker containers.

Suppose our BAM files and GTF files are stored in /yourdatafolder, and we're going to use rMATS-turbo to analyze them. Docker can't access these file for security reason. In order to make these file visible to Docker, we have to use option -v (<https://docs.docker.com/engine/tutorials/dockervolumes/#data-volumes>). This option will mount our local folder to docker's file system, and retrieve output from docker.

Note that, after mounting our folder to docker, docker can read this folder, and it can also write output file to this folder.

Examples

Suppose we have 4 samples in /yourdatafolder.

```
1 $ ls /yourdatafolder:
2 - b1.txt
3 - b2.txt
4 - 5.gtf
5 - 1.bam
6 - 2.bam
7 - 3.bam
8 - 4.bam
9
10 $ cat b1.txt:
11 /data/1.bam,/data/2.bam
12
13 $ cat b2.txt:
14 /data/3.bam,/data/4.bam
15
16 docker run -v /yourdatafolder:/data rmats:turbo01 --b1 /data/b1.txt \
17 --b2 /data/b2.txt --gtf /data/5.gtf --od /data/output -t paired \
18 --nthread 4 --readLength 101 --anchorLength 1
```

This command mounts the host directory, /yourdatafolder, into the container at /data. If the path /data already exists inside the container's image, the /yourdatafolder mount overlays but does not remove the pre-existing content. Once the mount is removed, the content is accessible again. This is consistent with the expected behavior of the mount command.

Accordingly, the absolute path of file should be adjusted. (e.g. b1.txt, 5.gtf, 2.bam, etc. changed to /data/b1.txt, /data/5.gtf, /data/2.bam, etc.)

Important note: The output folder /data/output will be written to /yourdatafolder/output.