

# RRBS Analysis Package for PBS

---

## Contents

|                                    |   |
|------------------------------------|---|
| 1 RRBS installation and usage..... | 1 |
| 2 Configuring RRBS.....            | 3 |
| 2.1 workflow.properties.....       | 3 |
| 2.2 mysql.properties.....          | 4 |
| 2.3 rrbs.properties.....           | 4 |
| 2.4 reference.xml.....             | 6 |
| 3 Contact and Support.....         | 8 |

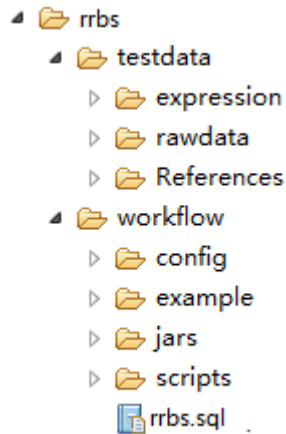
## 1 RRBS installation and usage

1. Install the required software and configure the environment. (see the prerequisites document)
2. Download the rrbs\_pipeline\_pbs.tar.gz from WBSA website (<http://wbsa.big.ac.cn/download/download.jsp>) put it under a given directory such as /home/test

**Note: In the following paragraph we will use “/home/test” to make examples.**

```
cd /home/test  
  
tar xzvf rrbs_pipeline_pbs.tar.gz
```

After uncompressed, you will see the directory structures like the following figure:



**Note:** The files in the “rrbs/workflow/example” use “/home/test” as relative directory.

| Folder name         | Description   |
|---------------------|---|
| testdata/expression | an expression file to test the “EXPRESSION” function  |
| testdata/rawdata    | <b>required</b> ,a source fastq file to test the package  |
| testdata/References | <b>required</b> ,a test reference for test the package which need to be configured in the reference.xml |
| workflow/config     | contains configuration template for the package   |
| workflow/example    | contains all the configuration files to test the package  |
| workflow/jars       | <b>required</b> , contains all the jars to run the package  |
| workflow/scripts    | <b>required</b> , contains all the scripts to run the package   |
| rrbs.sql            | <b>required</b> ,a database schema need to load into MySQL database                                     |

### 3. Import the MySQL script.

You need to login into Mysql database as root, by default the password is empty, so just click enter.

```
mysql -uroot -p
```

```
mysql> create database rrbs;

mysql>source /home/test/rrbs/workflow/rrbs.sql;

mysql> grant all on rrbs.* to 'pipeline'@'localhost' identified by 'pipeline'
```

### 4. Run the RRBS module

**Note: Before you run the following command, please ensure you have configured the files which are referred in section 2**

```
java -jar /home/test/rrbs/workflow/jars/workflow.jar \
/home/test/rrbs/workflow/example/rrbs.properties \
/home/test/rrbs/workflow/example/config
```

## 2 Configuring RRBS

The package uses four configuration files to store several parameters which needed by the RRBS pipeline; the configuration files list here from the “example” directory.

- workflow.properties
- mysql.properties
- rrbs.properties
- reference.xml

### 2.1 workflow.properties

**Note: the parameters are identified by bold font you need to change to your real directory. This is a PBS(Portable Batch System) package, so you need to configure the cluster.submit.queueName parameter otherwise the whole package will not run correctly.**

```
path.result=/home/test/rrbs/workflow/example/result
path.perl=/home/test/rrbs/workflow/scripts
path.reference=/home/test/rrbs/workflow/example/config/reference.xml
path.config=/home/test/rrbs/workflow/example/config

database.use=1
package.type=PBS
cluster.submit.queueName=workq
```

Table 1 workflow.properties

| Parameter      | Description                                | Example       |
|----------------|--|---------------|
| path.result    | Absolute path of the result directory      | result        |
| path.perl      | Absolute path of the perl script directory | scripts       |
| path.reference | Absolute path of the reference.xml file    | reference.xml |

|                          |  |        |
|--------------------------|--|--------|
| path.config              | Absolute path of directory of the configuration files                      | config |
| database.use             | Whether or not to use database. 1: yes, 0: no. Necessary in this workflow. | 1      |
| package.type             | Type of the cluster, we use PBS in this workflow.                          | PBS    |
| cluster.submit.queueName | Work queue name in PBS used to submit jobs.                                | workq  |

## 2.2 mysql.properties

**Note: If the MySQL database which you have installed is not located in your current machine, you need change "localhost" to the IP of that machine.**

```
jdbc.driverClassName=com.mysql.jdbc.Driver

jdbc.url=jdbc:mysql://localhost/rrbs?user=pipeline&password=pipeline&useUnicode=true&characterEncoding=UTF-8&failOverReadOnly=false
```

## 2.3 rrbs.properties

```
pipeline.type=rrbs
refld=d2
pValue=0.1
regalement=p1
tRichFile=/home/test/rrbs/testdata/rawdata/test_trich.fastq
aRichFile=
isTrimQV=1
isTrimAdaptor=1
minQualityValue=20
adaptorSeq=AGATCGGAAGAGC
minLength=35
minQualityValueForC=20
isFilterStartPoint=1
seedsLength=32
seedsMismatch=2
totalMismatch=4
lambdaFile=
expressionFile=/home/test/rrbs/testdata/expression/test.exp
teFile=
jar.package=/home/test/rrbs/workflow/jars/rrbs.jar
jar.picture=/home/test/rrbs/workflow/jars/pipelinepicture.jar
jar.html=/home/test/rrbs/workflow/jars/pipelinehtml.jar
```

The following table describes the parameters in the rrbs.properties

Table 2 parameters of rrbs.properties

| Parameter           | Description  | Example       |
|---------------------|--|---------------|
| pipeline.type       | The pipeline type of RRBS, <b>the value must be rrbs</b>   | rrbs          |
| refId               | Reference ID for the data, which is included in the file reference.xml.  | d2            |
| regElement          | Regulate element of the raw data including: p1 p2, p1 represents t-rich file,p2 represent a-rich file          | p1 p2         |
| tRichFile           | Absolute path of the t-rich file.  | trich.fastq   |
| aRichFile           | Absolute path of the a-rich file.  | arich.fastq   |
| isTrimQV            | Whether or not to filter the low quality bases of the raw data. 1:yes, 0: no.                                  | 1 0           |
| isTrimAdaptor       | Whether or not to filter the adaptor sequences   | 1 0           |
| minQualityValue     | The user could trim low quality bases from two ends if the base quality value is less than a threshold         | 20            |
| adaptorSeq          | Adaptor sequences  | AGATCGGAAGAGC |
| minLength           | If the read length is less than a preset value after above trimming and filtering, the read will be discarded. | 35            |
| minQualityValueForC | The minimum quality value is a threshold, less than which the C base is not considered to be a methylcytosine. | 20            |
| isFilterStartPoint  | Whether or not to remove the duplicated reads.   | 1 0           |
| lambdaFile          | Absolute path of lambda file.  | lambda.fa     |
| seedsLength         | BWA parameter: seeds length.   | 32            |
| seedsMismatch       | BWA parameter: maximum difference in the seed.   | 2             |
| totalMismatch       | BWA parameter: maximum differences in total reads.   | 4             |
| pValue              | p-value, if not choose to use lambda sequence to calculate it, users should offer it directly.                 | 0.005         |
| expressionFile      | A file with gene information and expression value.<br>#gene id #express value #chromosome                      |               |

|              |  |  |
|--------------|--|--|
|              | #gene start position (count from 0)<br>#gene end position (count from 1)<br>#strand  |  |
| teFile       | The upload TE data file format is as below:<br>#species #TE_id #chromosome<br>#strand #start position (count from 1)<br>#end position (count from 1) |  |
| jar.package  | for RRBS pipeline, <b>the value must be rrbs.jar</b>   |  |
| jar. picture | to generate pictures, <b>the value must be pipelinepicture.jar</b>   |  |
| jar.html     | to generate html result page, <b>the value must be pipelinehtml.jar</b>  |  |

## 2.4 reference.xml

**Note:** The parameter which identified by bold font you need to change to satisfy your real demands, for the left parameters and the whole XML format you should not change it in order to run the package correctly.

If you have multiple references, just add the <reference> element as the example file.

```
<?xml version="1.0" encoding="UTF-8"?>
<references>
<reference id="d2" name="Human">
  <desc>refence for demo data in RRBS,chrom 22</desc>
  <params>
    <param name="refDir">/home/test/rrbs/testdata/References/human/ref</param>
    <param
name="refC2TG2AFile">/home/test/rrbs/testdata/References/human/Ref_C-T_G-A/ref_all_C-T_G-A.fa</param>
    <param name="genelistChromDir">/home/test/rrbs/testdata/References/human/genes</param>
    <param name="repeatDir">/home/test/rrbs/testdata/References/human/repeats</param>
    <param
name="goNumberFile">/home/test/rrbs/testdata/References/human/GO/human_gene_GO.txt</param>
    <param
name="geneOntologyFile">/home/test/rrbs/testdata/References/human/GO/gene_ontology_ext.obo</param>
    <param name="refCGIDir">/home/test/rrbs/testdata/References/human/cpgislands</param>
    <param name="ideogramDir">/home/test/rrbs/testdata/References/human/ideogram</param>
  </params>
</reference>
```

</references>

Table 3 attributes of element <reference>

| Attribute | Description   |
|-----------|---|
| Id        | ID of the species, <b>The value must be unique, it will be used by the parameter "refid" in the rrbs.properties</b> |
| Name      | Name of the species.  |

Table 4 parameters of a reference < reference >

| Param            | Description  | Example   |
|------------------|--|---|
| refDir           | Absolute path of directory of raw reference files that must be standard .fa format and named chr*.fa.  | /home/test/rrbs/testdata/References/human/ref                               |
| refC2TG2AFile    | Absolute path of reference file that has merged all the raw reference files and converted C to T and G to A. The file must have been indexed by BWA before used.   | /home/test/rrbs/testdata/References/human/Ref_C-T_G-A/ref_all_C-T_G-A.fasta |
| genelistChromDir | Absolute path of gene file directory. The gene file is downloaded from UCSC. The genes on forward strand of each chromosome are stored in a file named chr*_C-T.gene. The genes on reverse strand of each chromosome are stored in a file named chr*_G-A.gene.<br>#geneid #geneid #chromosome #strand<br>#start #end #start #end #exon number<br>#first exon start pos #second exon start pos #first exon end pos #second exon end pos<br>All the start positions are count from 0 and the end positions are count from 1. | /home/test/rrbs/testdata/References/human/genes                             |
| repeatDir        | Absolute path of repeat file directory. It is downloaded from UCSC and in the standard format.   | /home/test/rrbs/testdata/References/human/repeats                           |
| goNumberFile     | File with GO information. Each line begins with a gene id followed by the GO number.   | /home/test/rrbs/testdata/References/human/GO/human_gene_GO.txt              |
| geneOntologyFile | Full ontology file, including  | /home/test/rrbs/testdata/References/human/GO/human_gene_GO.txt              |

|             |   |  |
|-------------|---|--|
|             | cross-products, inter-ontology links, and has part relationships. It could be downloaded at <a href="http://www.geneontology.org">http://www.geneontology.org</a> | stdata/References/human/GO/gene_ontology_ext.obo     |
| refCGIDir   | Absolute path of CpG islands file directory. It is downloaded from UCSC and in the standard format.   | /home/test/rrbs/testdata/References/human/cpgislands |
| ideogramDir | Absolute path of ideogram files of the species including: chromosomes.png, chromosomes.map.   | /home/test/rrbs/testdata/References/human/ideogram   |

### 3 Contact and Support

RRBS Analysis package is developed and maintained by [Beijing Institute of Genomics\(BIG\)](#), Chinese Academy of Sciences. If you have feedback or questions, please feel free to contact us at [wbsa@big.ac.cn](mailto:wbsa@big.ac.cn).