

# GUAVA Manual

Mayur Divate and Edwin Cheung

Contact: [mdivate@umac.mo](mailto:mdivate@umac.mo)

# Index

<b>1. How to install dependencies? .....</b>	<b>3</b>
1.1 Java 1.8 or latest .....	3
1.2 Bowtie version 1.1.2 .....	3
1.3 Python version 2.7 .....	4
1.4 MACS2 version 2.1.1.20160309 .....	4
1.5 SAMtools Version: 1.3.1 .....	4
1.6 R Version: >= 3.3.0 .....	5
<b>2. How to download and launch GUAVA? .....</b>	<b>5</b>
<b>3. GUAVA graphical user interface .....</b>	<b>6</b>
<b>4. GUAVA command-line interface .....</b>	<b>16</b>
4.1 Usage and option summary .....	16
<b>5. How to download genome fasta file? .....</b>	<b>17</b>
<b>6. How to create bowtie index of genome fasta file? .....</b>	<b>17</b>

# **GUAVA: A GUI tool for the Analysis and Visualization of ATAC-seq data**

GUAVA is a standalone GUI application for analyzing ATAC-seq data. GUAVA works on Linux and Mac OS. GUAVA is developed to help researcher with minimal or no Linux background for analyzing ATAC-seq data. This document contains all the information which is required to install dependencies and use GUAVA graphical and command line interface. This document also explains the GUAVA graphical user interface using published ATAC-seq data. It also provides procedure to create bowtie index from genome fasta for bioinformatics novice.

## **1. How to install dependencies?**

GUAVA depends on others tools in order to process ATAC-seq data e.g. bowtie for alignment. Those dependencies need to be installed on machine to use GUAVA. If any of the dependency is not found GUAVA will fail to start. Text followed by “\$” is a command. Once terminal (command line provided by OS) is launched, user can simply copy and paste on terminal those command to complete installation.

### **1.1 Java 1.8 or latest**

As GUAVA is developed in Java, it is required.

To innstall java on Mac OS,

- To Download java go to <https://java.com/en/download/>
- Double-click the pkg file to launch it
- Double-click on the package icon to launch install Wizard
- The Install Wizard displays the Welcome to Java installation screen. Click Next
- Click the Next button to continue the installation.
- Click Close to finish the installation process.

For more details, please follow this link

[https://www.java.com/en/download/help/mac\\_install.xml](https://www.java.com/en/download/help/mac_install.xml)

To install Java on Linux OS

copy paste following command to the terminal

```
$ sudo apt-get install oracle-java8-installer
```

Or follow the link:

[https://java.com/en/download/help/linux\\_x64\\_install.xml](https://java.com/en/download/help/linux_x64_install.xml)

### **1.2 Bowtie version 1.1.2**

- Download bowtie from here  
<https://sourceforge.net/projects/bowtie-bio/files/bowtie/1.1.2/>  
Linux OS : bowtie-1.1.2-linux-x86\_64.zip  
Mac OS : bowtie-1.1.2-macos-x86\_64.zip
- Copy downloaded bowtie file or the file path i.e. location and to paste it on terminal use,  
Mac => command + v  
Linux => ctrl + shift + v
- Launch the terminal and Use following commands to terminal to install bowtie

```
$ cp <bowtie file path> ~/
$ cd ~/
$ unzip bowtie-1.1.2*.zip
$ cd bowtie-1.1.2/
```

Mac OS

```
$ echo "export PATH=\$PATH:"`pwd` | cat - >>
~/.bash_profile
$ source ~/.bash_profile
```

Linux OS

```
$ echo "export PATH=\$PATH:"`pwd` | cat - >> ~/.bashrc
$ source ~/.bashrc
```

### 1.3 Python version 2.7

Python is required for macs2 installation.

Mac OS

- Download the Mac OS X 64-bit/32-bit installer (not the PPC installer) from the Python website  
<https://www.python.org/downloads/release/python-2711/>.
- Double-click the python-2.7.11-macosx10.6.pkg file in the Downloads folder.
- If you have Gatekeeper enabled, the installation will be blocked. Open System Preferences > Security & Privacy and click Open Anyway.
- Click Continue, Agree and Install buttons in the Install Python window.

Linux OS, Use following command to install python

```
$ sudo apt-get install python
```

### 1.4 MACS2 version 2.1.1.20160309

To install MACS2 on Mac or Linux OS use command below

```
$ pip install --user MACS2
```

### 1.5 SAMtools Version: 1.3.1

- Download samtools-1.3.1.tar.bz2 link:

<https://sourceforge.net/projects/samtools/files/samtools/1.3.1/>

- Copy downloaded SAMtools file or copy the samtools-1.3.1.tar.bz2 file path  
To paste path on terminal use  
Mac => command + v  
Linux => ctrl + shift + v
- Open the terminal
- Use copy following commands to terminal and hit enter  
\$ cp <samtools file path> ~/  
\$ cd ~/  
\$ tar jxvf samtools-1.3.1.tar.bz2  
\$ cd samtools-1.3.1  
\$ make
- Mac  
\$ echo "export PATH=\\$PATH:"`pwd` | cat - >> ~/.bash\_profile  
\$ ~/.bash\_profile
- Linux  
\$ echo "export PATH=\\$PATH:"`pwd` | cat - >> ~/.bashrc  
\$ source ~/.bashrc

## 1.6 R Version: >= 3.3.0

Mac user can follow the link below for video tutorial to install R. If R installation is found on system, GUAVA will install bioconductor packages automatically.

- [https://youtu.be/cX532N\\_XLIs?list=PLqzoL9-eJTNBDdKgJgJzaQcY60XmsXAHU](https://youtu.be/cX532N_XLIs?list=PLqzoL9-eJTNBDdKgJgJzaQcY60XmsXAHU)

Linux user can use following link for installation

- <https://cran.r-project.org>

## 2. How to download and launch GUAVA?

To download GUAVA use from GitHub

<https://github.com/MayurDivate/GUAVA>

There is no need to install GUAVA it can be easily launched as described below.

Unzip downloaded GUAVA package

```
$ cd <GUAVA folder>
```

To launch GUAVA graphical user interface use following command

```
$ cd <GUAVA folder>
```

```
$ java -jar GUAVA.jar
```

### 3. GUAVA graphical user interface

We used SRR891275[1] dataset to demonstrate GUAVA graphical user interface.

The image shows the GUAVA graphical user interface. It has a title bar with 'GUAVA' and standard window controls. Below the title bar is a menu bar with 'Help' and 'About'. The main area is divided into several sections:

- Input fastq reads:** Contains two buttons, 'R1 fastq' and 'R2 fastq', each followed by a text input field. The R1 field contains '/path/R1.fastq' and the R2 field contains '/path/R2.fastq'.
- Adapter Trimming:** Contains a checked checkbox 'Trim adapter ?'. To its right are three input fields: 'Maximum Ns' with value '2', 'Minimum Read length' with value '30', and 'Error Rate' with value '0.1'. Below this is an unchecked checkbox 'Nextera XT Adapter Sequence' followed by 'OR' and 'Adapter Sequence' with an empty text field.
- Alignment Parameters:** Contains a button 'Bowtie v1 Genome Index' followed by a text field with '/path/bowtie1Index.ebwt'. Below this are three input fields: 'Maximum insert size' with value '2000', 'No. of genomic hits (m)' with a spinner set to '1', and 'Genome Assembly' with a dropdown menu showing '-select-'.
- Chromosome Filtering:** Contains two checkboxes: 'chrM (recommended)' (checked) and 'chromosome Y' (unchecked). To their right are two input fields: 'RAM ( GB )' with a spinner set to '1' and 'CPU units' with a spinner set to '1'.
- Peak Calling and other parameters:** Contains a dropdown menu 'q value' set to '0.05' and a text field 'Organism :'. Below these is a button 'Output Folder' followed by a text field containing '/path/output\_dir'.

At the bottom of the window are two large buttons: 'Start Analysis' and 'Reset All'.

**Figure 1: Graphical user interface of GUAVA**

User can input the FASTQ files and bowtie index of genome fasta file using the buttons for R1 fastq, R2 fastq and Bowtie v1 Genome Index, respectively. If user select “Trim Data?” check box, then adapter will be trimmed from sequencing reads using cutadapt(Martin, 2011) before aligning them to genome. User can either select Nextera XT adapter or can provide custom adapter sequence. Reads with more “Maximum Ns” and less than “Minimum Read length” after adapter trimming will be filtered. Maximum distance for mapping of read mates can be set using “Maximum insert size” input field. To exclude reads pairs which are exceeding desired number of alignments (m) limits can be set

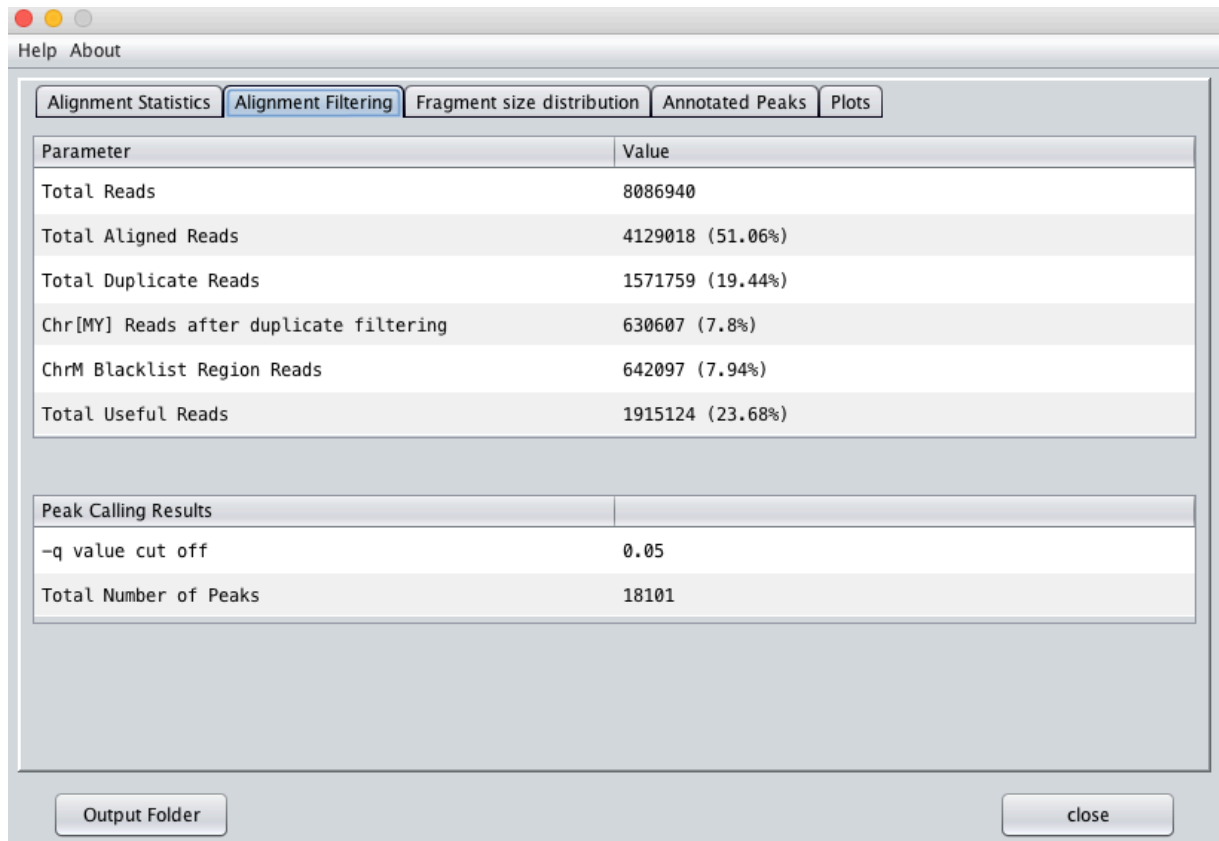
using “No. of genomic hits(m)” input field. User can specify the assembly or version name using genome assembly drop down menu. If chromosome checkbox is selected, reads mapping to that chromosome will be discarded. Dropdown menu is provided to select p or q value and input field next it can be used to set cutoff to filter MACS2(Feng, et al., 2012) peaks. RAM and threads can be set using box 5 and 6 respectively. Start analysis button will start ATAC-seq data processing, if provided inputs values are valid else invalid values will be highlighted with red color.

Parameter	Value
R1 fastq	SRR891275_1.fastq
R2 Fastq	SRR891275_2.fastq
Genome	hg19.fa
Maximum Insert Size	2000
Maximum genomic hits	1
Total Reads	8086940
Total Aligned Reads	4129018 (51.06%)
Total Reads Failed to Align	1819172 (22.5%)
Total Suppressed Reads	2138750 (26.45%)
Total chrM Reads	4344566

**Figure 2: Result interface tab for alignment**

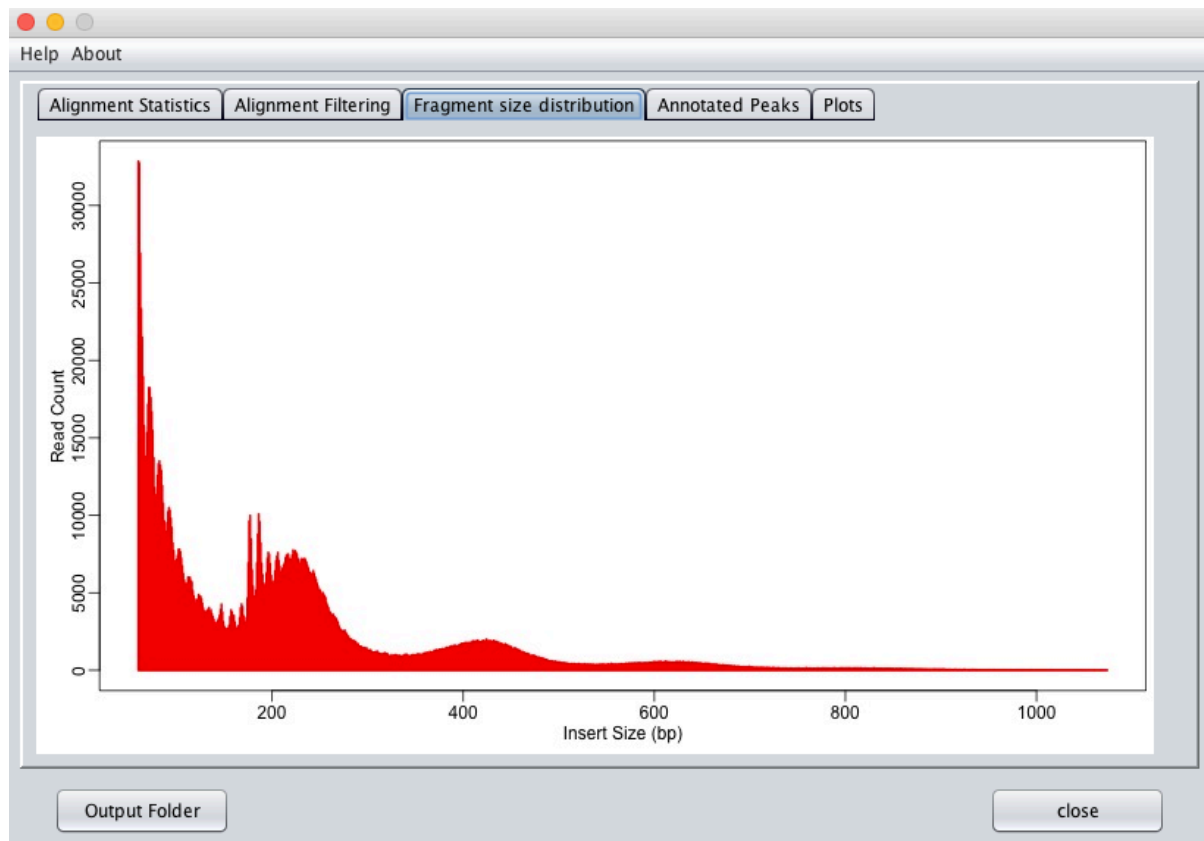
This includes the input summary: names of FASTQ files and genome assembly used for analysis, maximum insert size and genomic hits. Total number of reads, aligned reads un aligned reads , fail to align due to m limit. Reads mapping to chr M and If chrY selected.





**Figure 3: Alignment filtering and peak calling summary**

The first table provides the number and percentage of reads relative to the total input reads i) aligned reads ii) duplicate reads iii) reads aligned to chromosome M and Y (based on selection) iv) reads mapped in blacklisted regions v) Useful reads represents the actual reads that qualifies for further processing. Second table gives you total number of peaks called and p/q value cutoff as per input



**Figure 4: Fragment size distribution graph**

This graph gives fragment size distribution calculated using the Picard(Institute) tool. The x-axis shows the size of fragment in bp and while the y-axis shows the read count.

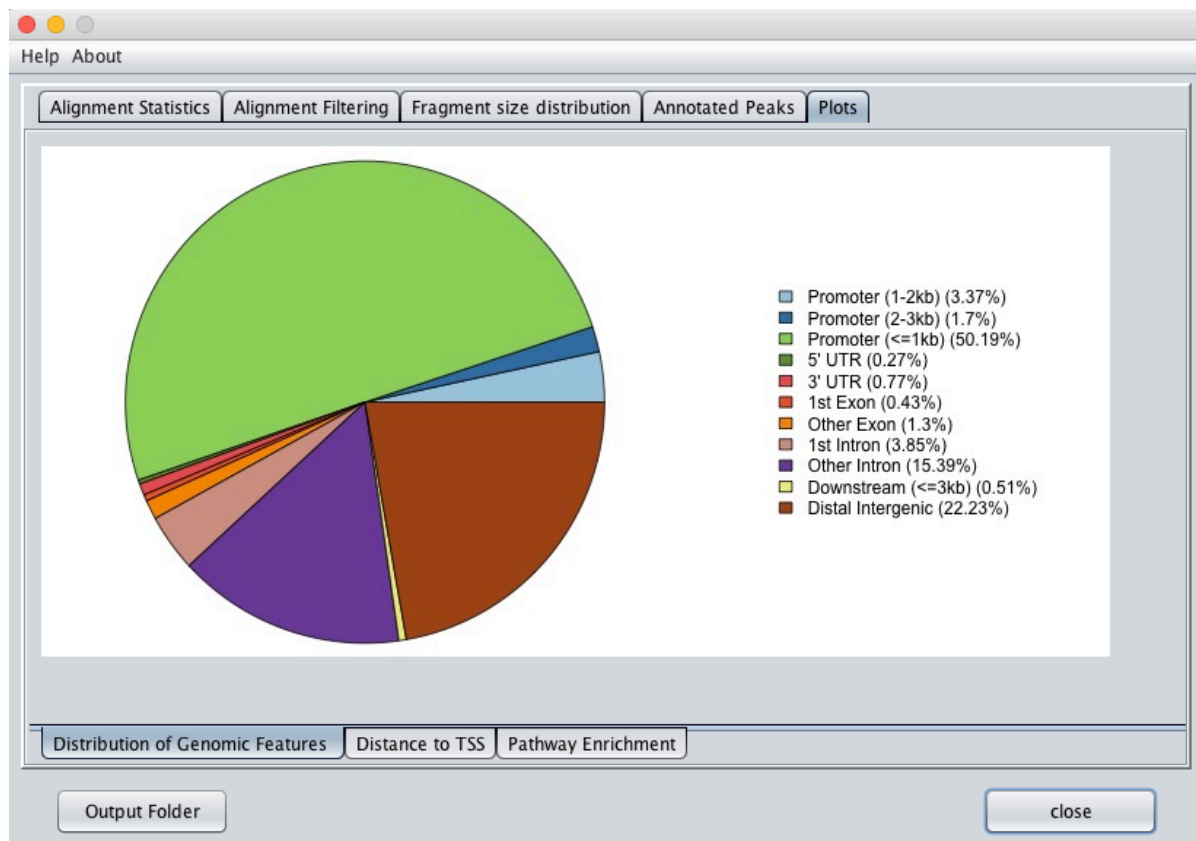
Chromo...	Start	End	Length	Pileup He...	-log10(p)	-log10(q)	Annotation	Distance...	Gene Sy...
chr4	1699305...	1699320...	1454	52	114.643	109.196	Promoter...	0	CBR4
chr5	1724830...	1724834...	423	15	26.949	23.154	Promoter...	0	CREBRF
chr5	1724838...	1724843...	576	12	20.877	17.272	Promoter...	-10	CREBRF
chr5	34915396	34916104	709	37	77.008	72.16	Promoter...	0	BRX1
chr6	90031465	90031663	199	4	6.371	3.48	Distal Int...	-6447	GABRR2
chr6	1417857...	1417859...	199	5	7.999	4.984	Distal Int...	624024	NMBR
chr6	1418754...	1418756...	199	4	6.371	3.48	Distal Int...	534290	NMBR
chr6	90008637	90008835	199	5	7.999	4.984	Intron (uc...	16183	GABRR2
chr6	1420290...	1420292...	199	3	4.822	2.082	Intron (uc...	380703	NMBR
chr6	1421962...	1421964...	199	4	6.371	3.48	Intron (uc...	213444	NMBR
chr6	1423104...	1423106...	199	4	6.371	3.48	Intron (uc...	99258	NMBR
chr6	1393501...	1393504...	244	8	13.248	9.944	Promoter...	362	ABRACL
chr6	36164879	36165154	276	6	9.693	6.572	Promoter...	329	BRPF3
chr6	1418059...	1418062...	286	5	7.999	4.984	Distal Int...	603711	NMBR
chr6	1418839...	1418842...	286	5	7.999	4.984	Distal Int...	525664	NMBR
chr6	1418047...	1418051...	463	44	94.304	89.168	Distal Int...	604758	NMBR
chr6	70505888	70506356	469	11	18.915	15.377	Promoter...	178	LMBRD1
chr6	1393496...	1393501...	495	11	18.915	15.377	Promoter...	0	ABRACL
chr6	1421695...	1421700...	528	9	15.097	11.71	Intron (uc...	239889	NMBR
chr6	32939316	32939921	606	6	9.693	6.572	Promoter...	0	BRD2
chr6	42531463	42532076	614	8	13.248	9.944	Promoter...	0	UBR2
chr6	70506558	70507234	677	60	135.502	129.738	Promoter...	0	LMBRD1
chr7	45151254	45151480	227	8	13.248	9.944	Promoter...	0	TBRG4
chr7	2595321	2595623	303	5	7.999	4.984	Promoter...	0	BRAT1
chr7	1566856...	1566864...	763	12	20.877	17.272	Promoter...	0	LMBR1
chr7	07010405	07011180	776	7	11.445	8.228	Promoter...	0	BRD2

**Figure 5: Peak annotation table**

Annotated peaks can be browsed using peak annotation table. Gene name search boxed helps to filter peaks using gene symbol column. Table can be sorted by clicking on header of column. View in IGV button load the data tracks in IGV(Robinson, et al., 2011) and sets IGV location same as peak coordinates.

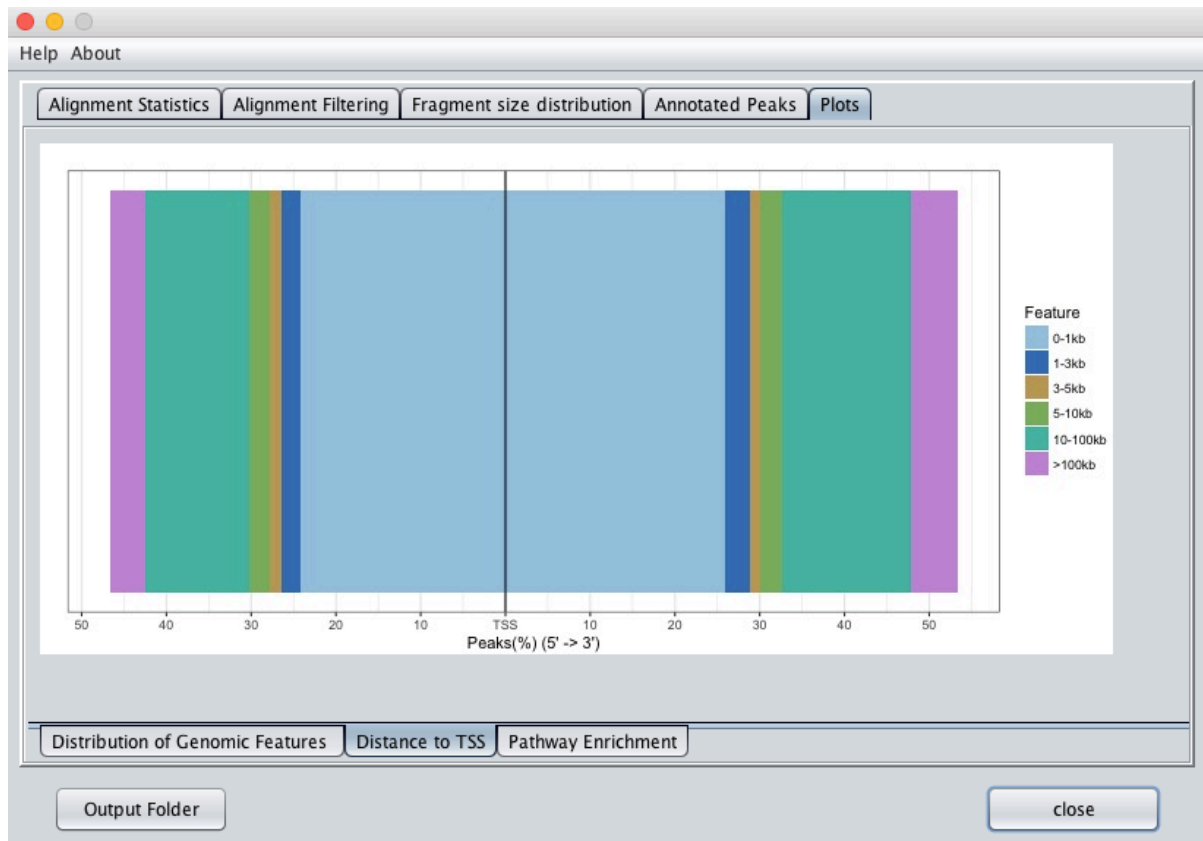


Figure 6: Peak visualization using IGV



**Figure 7: Peak annotation pie chart**

Pie chart showing genomic feature wise peak distribution obtained using ChIPseeker bioconductor package. In this example, more than 50% of the peaks are in the promoter region. The promoter region is divided into three sub-categories: 0 to 1 kb, 1 to 2 kb and 2 to 3 kb from the TSS.



**Figure 8: Peak distribution based on relative distance to the TSS**

Graph showing the percentage of peaks in a given range from the TSS. Different colors indicate the different range for distance and y-axis shows the percentage of total peaks. From above graph, around 50% peaks are within 0-1 kb (sky blue color) from the TSS.



**Figure 9: Pathway enrichment analysis**

This graph shows most significantly enriched pathway in given sample using ReactomePA(Yu and He, 2016) bioconductor package. Here, Cell cycle pathway enriched with  $1 \times 10^{-6}$  adjusted p-value.

## 4. GUAVA command-line interface

Like graphical user interface also has command line interface. Command line user interface make GUAVA flexible and therefore it can be easily integrated into existing pipelines. Also provides flexibility for running GUAVA through resource manger or job scheduler systems such as SLURM.

To print help message

```
$ java -jar GUAVA.jar -h
```

### 4.1 Usage and option summary

Usage: \$ java -jar **ATAC\_GUI.jar** [options]\*

Options	Description
<b>R1</b>	Path to the FASTQ file containing upstream mates
<b>R2</b>	Path to the FASTQ file containing upstream mates
<b>g</b>	Path to bowtie index of genome fasta file
<b>a</b>	Genome assembly version [hg18,hg18,hg38,mm9,mm10]
value	p   q value for MACS2 peak filtering default: q
c   cutoff	Cutoff for p/q value e.g. 0.05, 5E-2 default: 0.05
X	Maximum distance from each other at which read mates can map to the genome default: 2000
m	Report alignment for pair, if maximum number of reportable alignments for pair is less or equal to m default: 1
O   outdir	Path to the output directory default: current directory
ram	RAM memory to use in GBs default: 1
cpu	Number of threads to use default: 1
chrM	Remove(T) or keep(F) reads mapping to mitochondrial chromosome default: T
chrY	Remove(T) or keep(F) reads mapping to chromosome Y default: F
H   help	Print help mesage

Table 1: Usage and options for GUAVA command line interface. Compulsory options are shown in blue color.



## 5. How to download genome fasta file?

To download genome fasta files please follow the following links.

Human: <http://hgdownload.soe.ucsc.edu/downloads.html#human>

Mouse: <http://hgdownload.soe.ucsc.edu/downloads.html#mouse>

Select genome assembly that you want to download, then click on “Full data set” and download \*.fa.gz or chromFa.tar.gz (one chromosome per file) file.

To decompress file use the following command:

```
$ gzip -d <file name>
```

If it is one chromosome per file then make one single fasta file all chromosomes using cat command

e.g.

```
$ cat file1.fa file2.fa > genomeName.fasta
```

## 6. How to create bowtie index of genome fasta file?

If you already have genome fasta file, please follow the command below to create bowtie genome index. Bowtie uses this index to speed up alignment process.

```
$ cd <path to genome fasta file>
```

```
$ bowtie-build <genome.fasta> <genome>
```

**Note:** This is a time consuming step

## References

Feng, J., *et al.* Identifying ChIP-seq enrichment using MACS. *Nat Protoc* 2012;7(9):1728-1740.

Institute, B. Picard. <http://broadinstitute.github.io/picard/>.

Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnetjournal*, North America 2011.

Robinson, J.T., *et al.* Integrative genomics viewer. *Nat Biotechnol* 2011;29(1):24-26.

Yu, G. and He, Q.Y. ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. *Mol Biosyst* 2016;12(2):477-479.