

GUAVA Manual

Mayur Divate and Edwin Cheung

Index

Download	4
Getting help and reporting issues	4
How to open Terminal?	5
Install	5
Install dependencies	6
#1 Install R	6
#2 Install other dependencies and R packages	6
#3 Install MACS2	6
How to start GUAVA?	7
Graphical user interface of GUAVA	8
ATAC-seq data analysis program: Parameters	8
Output interface for GUAVA ATAC-seq data analysis	11
ATAC-seq differential analysis program: parameters	15
Output interface of GUAVA ATAC-seq differential analysis	17
How to download genome fasta file?	20
How to create a bowtie index of genome fasta file?	20

GUAVA: a GUI tool for the Analysis and Visualization of ATAC-seq data

In nutshell, GUAVA is a standalone GUI tool for processing, analyzing and visualizing ATAC-seq data. A user can start GUAVA analysis with raw reads to identify ATAC-seq signals. Then ATAC-seq signals from two conditions can be compared using GUAVA to identify genomic loci with differentially enriched ATAC-seq signals. Furthermore, GUAVA also provides gene ontology and pathways analysis. Since to use GUAVA requires only several clicks and no learning curve, it will help novice bioinformatics researchers and biologist with minimal computer skills to analyze ATAC-seq data. Therefore, we believe that GUAVA is a powerful and time saving tool for ATAC-seq data analysis. GUAVA setup contains a script to configure and install dependencies which facilitates the GUAVA installation. GUAVA works on Linux and Mac OS.

This document contains all the information that is required to install and use GUAVA.

GUAVA is developed in the Edwin's laboratory at University of Macau.

Download

Step 1: Go to the link: <https://github.com/MayurDivate/GUAVA/releases>

Step 2: Click on the `Source Code (zip)`

Step 3: This will save GUAVA zip package in the downloads folder

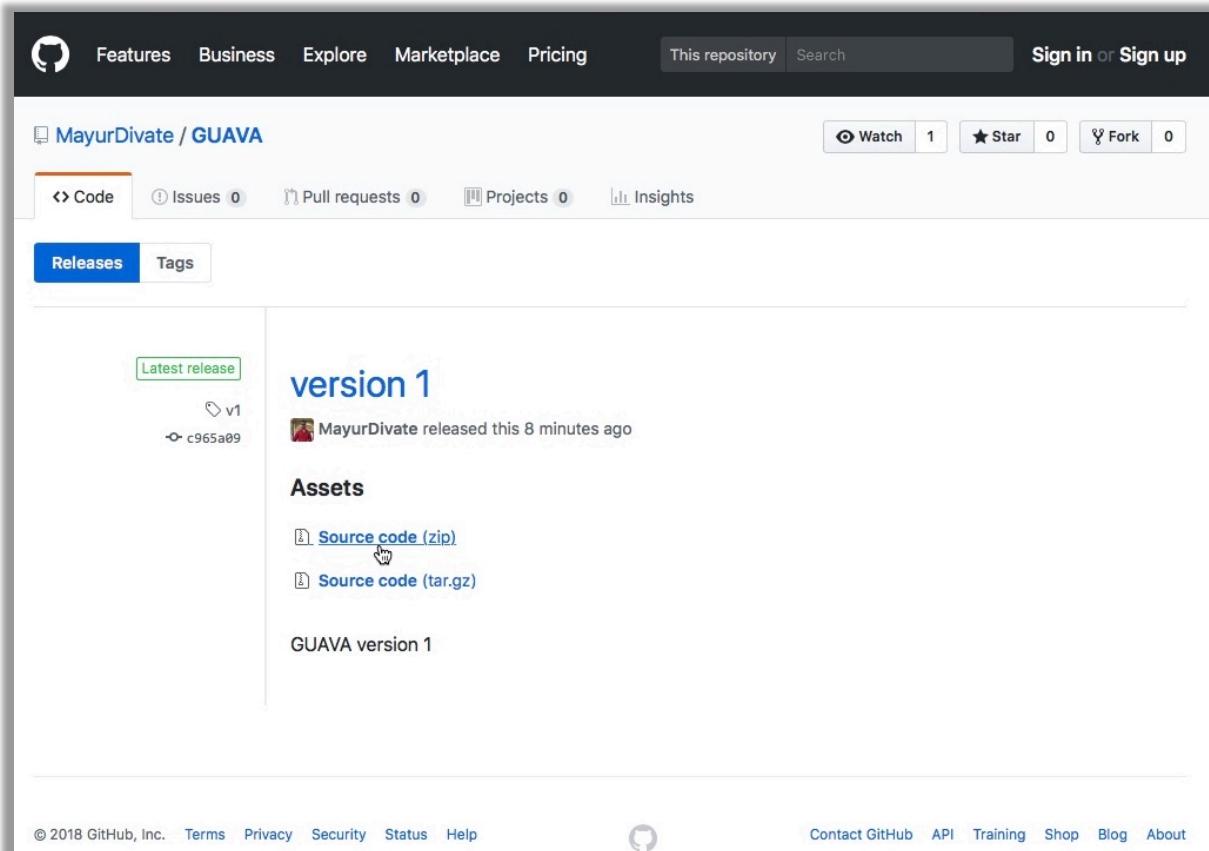


Figure1: GUAVA - GitHub project page

Go to GUAVA project page on GitHub and click on the ‘releases’ option. After that, click on the ‘Source code (zip)’ under the latest release.

If want download source code, follow the link below.

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

Getting help and reporting issues

If you face difficulties in installing or using GUAVA please report the issue here

<https://github.com/MayurDivate/GUAVASourceCode/issues>

How to open Terminal?

After downloading GUAVA, the user will need to open terminal to install terminal. Here we describe procedure to open terminal for the non-bioinformatics users.

MAC

1. Open the Finder
2. Click on the `Go` in the menu bar then select the `Utilities` folder
3. Double-click on the Terminal icon

Linux

1. Press windows key on keyboard
 2. Type “terminal” in the search box
 3. Click on the Terminal icon
- OR press “**Ctrl + Alt + T**” this will open then terminal

Install

Let's us assume that the user wants to install GUAVA in the home folder. Therefore, before proceeding, first copy / move downloaded GUAVA package to the home folder. If the downloaded package is in the folder `Downloads`, then type the following commands to unzip package on the terminal.

```
mv ~/Downloads/GUAVA-1.zip ~/  
cd ~/  
unzip GUAVA-1.zip  
mv GUAVA-v1 GUAVA
```

NOTE: If the downloaded GUAVA package is in the different folder than ‘Downloads’, you will have to use complete path of that folder instead of ~/Downloads/GUAVA-master.zip. To copy path, simply copy the downloaded package and paste it on the terminal.

Install dependencies

GUAVA depends on other tools in order to process ATAC-seq data (e.g. bowtie for alignment). If any of the dependency is not found on system, GUAVA will not work properly. Therefore, to help users to install dependencies, we have written a program (configure.sh) which automatically downloads and installs dependencies. However, the user need to install R and MACS2 manually due to the technical reasons.

Step 1: Install R

MAC

1. Download R follow this link => <https://cran.r-project.org/bin/macosx/>
2. Click on the R-X.X.X.pkg file link (e.g. R-3.4.3.pkg)
3. Double click on the downloaded file and follow the instructions

Linux

1. Open the terminal
2. Type command ` sudo apt-get install r-base ` and press enter

To know more about R, follow the link <https://cran.r-project.org/bin/linux/>. Choose appropriate Linux OS type

Step 2: Install other dependencies

To run configure.sh use the following commands on the terminal.

```
cd ~/GUAVA  
sh ./configure.sh
```

Note: This may take a while to finish. Also, you will need to press 'enter' several times to continue. Additionally, answer all question with 'yes'

Step 3: Install MACS2

Now, to install last dependency MACS2 use following commands on the terminal

```
cd ~/GUAVA  
python get-pip.py  
pip install MACS2
```

Error: Sometime MACS2 fails to install **Numpy**. In such a situation run `pip install numpy` first then try to install MACS2.

NOTE: If **permission denied**, type 'sudo' at the beginning of the commands. Then, in order to proceed you have enter your password

That is the end of installation part. Now, GUAVA is ready to process ATAC-seq data.

How to start GUAVA?

After successful installation of GUAVA, user can start GUAVA using following command on the terminal. This will open home window of GUAVA where you can choose program to proceed further.

```
cd ~/GUAVA  
java -jar GUAVA.jar
```

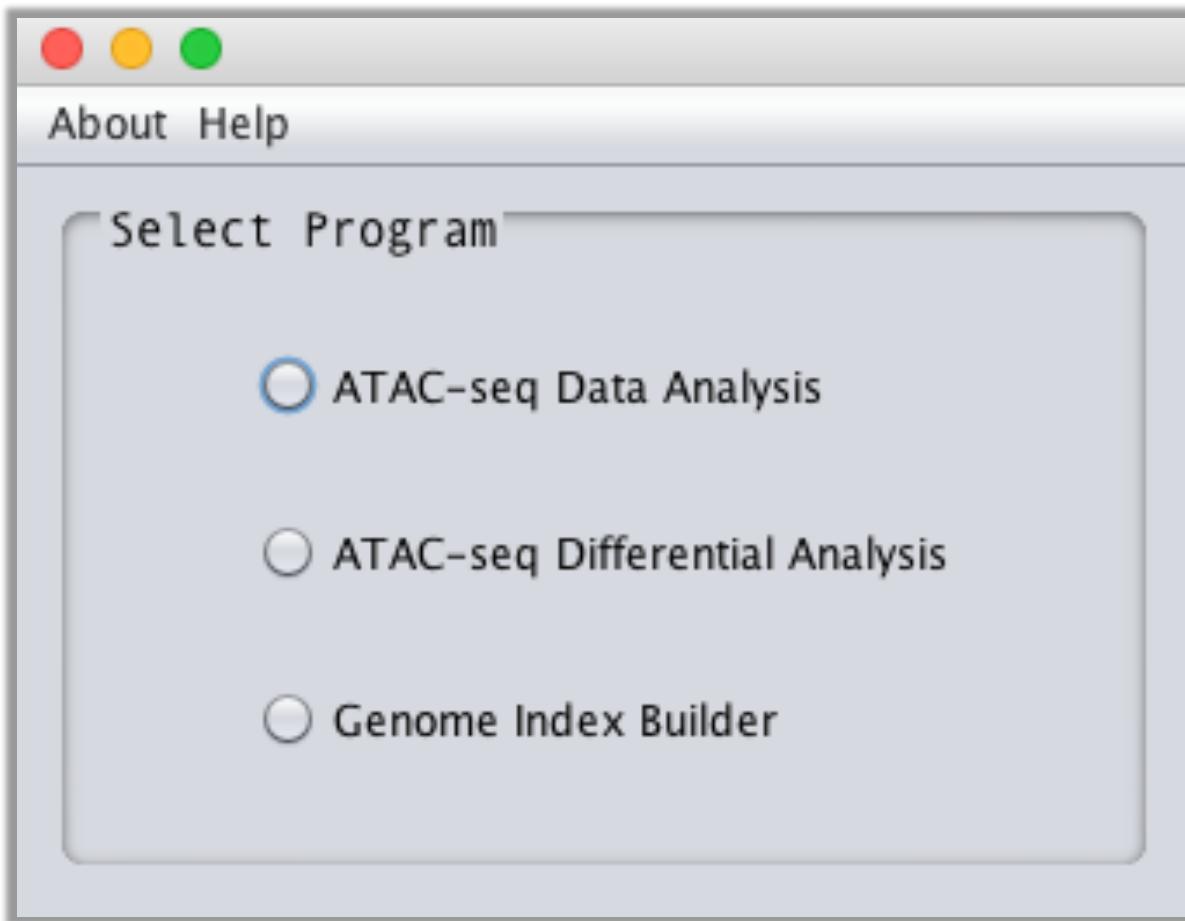


Figure 2: GUAVA – home window

Choose ATAC-seq Data Analysis to process ATAC-seq data, ATAC-seq Differential Analysis to compare ATAC-seq datasets or Genome Index Builder to generate genome index.

Graphical user interface of GUAVA

GUAVA tool has three programs

- 1) ATAC-seq data analysis: to process raw ATAC-seq sequencing reads
- 2) ATAC-seq differential analysis: to compare ATAC-seq signals.
- 3) Genome index builder

When GUAVA GUI is evoked it open GUAVA home window (Figure 2). Then user have to choose one of the above programs to proceed further. Finally, based on the selection of program, the desire input window will be opened (Figure 3, 13 and 22).

ATAC-seq data analysis program

This program accepts raw ATAC-seq reads as an input. Before aligning reads to genome, it trims adapter sequence from reads using cutadapt only if trimming option is selected. After that it filters unsuitable reads for ATAC-seq analysis such as duplicate reads. Additionally, the user can exclude certain chromosomes from the analysis using ‘show chromosomes’ button (Figure 4). Next, it uses MACS2 to identify ATAC-seq peaks. Finally, it performs functional annotation on the ATAC-seq peaks.

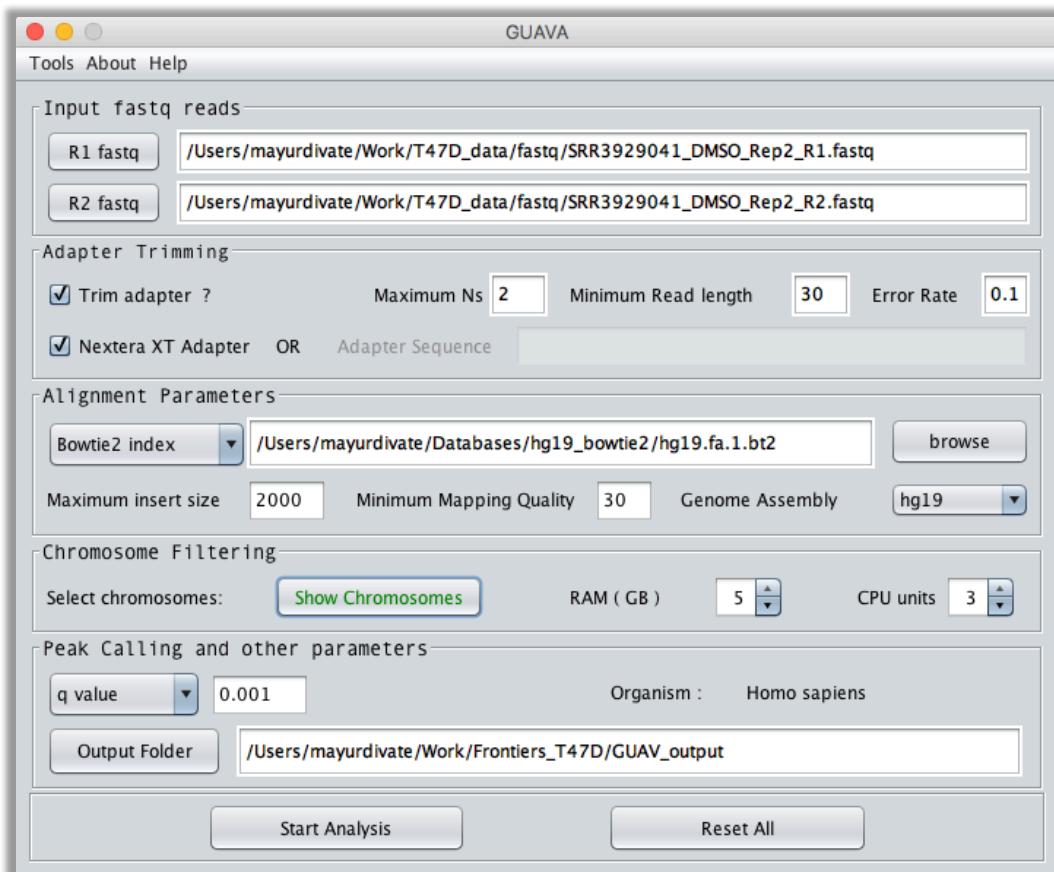


Figure 3. Input form of GUAVA ATAC-seq data analysis program

Input window of ATAC-seq data analysis program to upload input files such as fastq, genome index and set parameters such as insert size, p/q value etc.

R1 fastq and R2 fastq: Buttons to upload read1 and read2 fastq files of ATAC-seq data
Trim adapter: check this option if reads contains adapter
Maximum Ns: If any read contains more than specified number Ns after the adapter trimming, that read pair will be discarded (default 2)
Minimum read length: If any read is shorter than specified length after adapter trimming, that read pair will be discarded (default 30)
Error Rate: Allowed number of mismatches as a fraction of length. For example, if error rate is 0.1 then 1 mismatch is allowed for 10bp match of adapter sequence (default 0.1)
Nextera XT adapter: User can select this option if adapter used for ATAC-seq is a Nextera XT adapter (default adapter)
Adapter sequence: Option to specify custom adapter sequence when Nextera XT adapter is not used for library preparation.
Bowtie V1 or Bowtie V2 index: If you want to use bowtie for read mapping select “Bowtie index” from drop down menu else select “Bowtie2 index” to use bowtie2. Then using ‘browse’ button upload the appropriate genome index file (bowtie or bowtie2 index). Please see section ‘how to create genome index’ to know more about genome index and genome index builder tool.
Maximum insert size: Maximum insert size in base pair allowed for paired end alignment (default 2000)
Maximum genomic hits or Mapping quality: Maximum genomic hit (bowtie) and Minimum Mapping quality (bowtie2) to discard reads pairs which has multiple alignments. default Maximum genomic hits =1 and Mapping quality = 30
Genome assembly: Select the correct genome build from drop down menu e.g. hg19 and same build will be used for peak annotation and functional analysis.
Show chromosomes: Button to exclude reads mapping to specific chromosomes such mitochondrial chromosome. After clicking this button, it will open a new window (Figure 4) there user can select desired chromosomes.
RAM: RAM in GB to be used by GUAVA (default 1)
CPU units: Number of CPU units to be used by GUAVA (default 1)
p or q value: Select appropriate value from drop down menu and specify the cut off value in box next to it. This will be used by MACS2 to filter peaks (default q value)
Output folder: Select folder to save GUAVA ATAC-seq data analysis results
Reset All: Button to set all parameters to default value.
Start Analysis: click this button to start ‘ATAC-seq data analysis’ program. If all provided options are valid then GUAVA will start analysis.

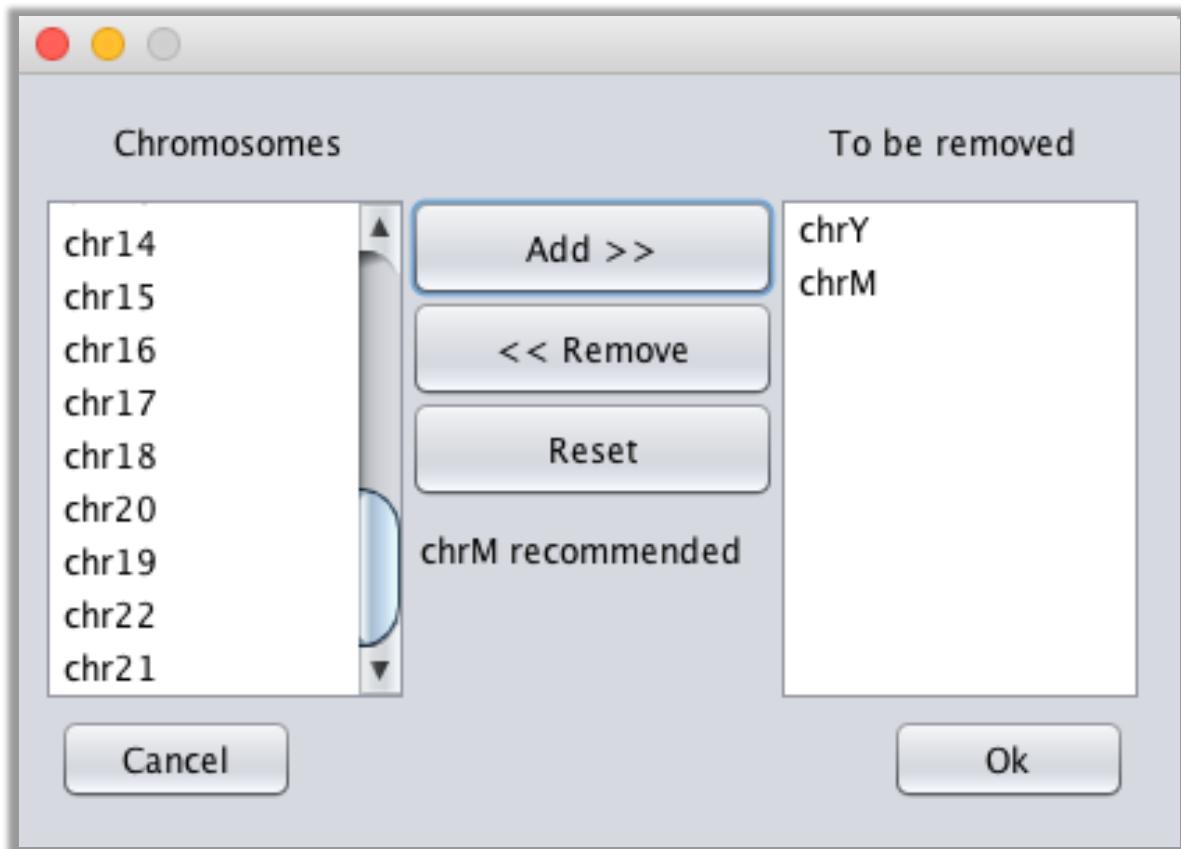


Figure 4. interface to select chromosomes for alignment filtering.

Here user can add desired chromosomes to the ‘to be removed’ list to discard reads aligning to those chromosomes.

Output interface for ATAC-seq data analysis program

Once GUAVA finishes analysis it shows results on tabular output interface (Figure 5-12). Also facilitates the visualization of ATAC-seq signal on IGV browser.

Parameter	Value
R1 fastq	SRR3929041_DMSO_Rep2_R1.fastq
R2 fastq	SRR3929041_DMSO_Rep2_R2.fastq
Genome index	hg19.fa
Maximum Insert Size	2000
Minimum Mapping Quality	30
Total Reads	94208170
Total Aligned Reads	75771166 (80.43%)
Total Reads Failed to Align	7719282 (8.19%)
Total Reads with Low Mapping Quality	10717722 (11.38%)
Total chrY Reads	103281 (0.11%)
Total chrM Reads	2907706 (3.09%)

Output Folder close

Figure 5: Input summary and alignment statistics This tab provides reads mapping statistics (e.g. total number of reads mapped to genome along) with summary of input files and parameters

Parameter	Value
Total Reads	94208170
Total Aligned Reads	75771166 (80.43%)
Total Duplicate Reads	13855635 (14.71%)
Chr* Reads after duplicate filtering	823728 (0.87%)
Blacklist Region Reads	134548 (0.14%)
Total Useful Reads	60956787 (64.7%)

Peak Calling Results	
-q value cut off	0.001
Total Number of Peaks	147675

Output Folder close

Figure 6: Read filtering and peak calling summary. This tab has two tables 1) to provide alignment filtering statistics e.g. duplicate, useful reads etc. 2) to provide summary of MACS2 peak calling

About Help

Alignment Statistics		Alignment Filtering		Fragment Size Distribution		Annotated Peaks		Plot	Gene Ontologies	Pathways
Chr	Start	End	Length	Pileup H...	-log10(p)	-log10(q)	Annotation	Distance...	Gene Symbol	
chr16	56964546	56964757	212	21	8.515	3.989	upstream	1244	HERPUD1	
chr16	56965543	56967221	1679	318	454.93	57.839	overlapStart	0	HERPUD1	
chr16	69419571	69420171	601	140	152.914	25.565	overlapStart	0	TERF2	
chr16	75681150	75682850	1701	197	242.62	35.9	overlapStart	0	TERF2IP	
chr17	5389073	5390365	1293	145	160.433	26.472	overlapStart	0	DERL2	
chr17	80593318	8060381	1064	144	158.923	26.29	overlapStart	0	PER1	
chr17	8062038	8062418	381	30	15.515	5.621	upstream	2359	PER1	
chr17	27181520	27182297	778	181	216.617	32.999	overlapStart	0	ERAL1	
chr17	37843890	37844559	670	128	135.188	23.389	overlapStart	0	ERBB2	
chr17	62206619	62208372	1754	200	247.559	36.444	overlapStart	0	ERN1	
chr17	62210911	62211316	406	85	76.019	15.593	upstream	3408	ERN1	
chr18	44702457	44703030	574	91	83.803	16.681	overlapStart	0	IER3IP1	
chr18	61220179	61220752	574	130	138.11	23.752	upstream	2640	SERPINB12	
chr18	61369129	61369337	209	22	9.226	4.17	upstream	747	SERPINB11	
chr18	61549751	61550433	683	128	135.188	23.389	upstream	4505	SERPINB2	
chr18	61624059	61624269	211	25	11.466	4.714	upstream	3204	SERPINB8	
chr19	344710	345099	390	36	20.822	6.709	overlapStart	0	MIER2	
chr19	348191	348887	697	51	35.785	9.428	upstream	3399	MIER2	
chr19	5718514	5719073	560	43	27.533	7.978	upstream	1614	CATSPERD	
chr19	5719905	5720960	1056	132	141.045	24.115	overlapStart	0	CATSPERD	
chr19	8274038	8274713	676	41	25.564	7.615	overlapStart	0	CERS4	
chr19	10443312	10444695	1384	94	87.759	17.225	overlapStart	0	RAVER1	
chr19	10445261	10446000	740	115	116.531	21.032	upstream	946	RAVER1	
chr19	13260496	13261472	977	90	82.493	16.5	overlapStart	0	IER2	
chr19	14585683	14586821	1139	36	20.822	6.709	overlapStart	0	PTGER1	
chr19	14587525	14587753	229	26	12.246	4.895	upstream	1350	PTGER1	
chr19	14588106	14588316	211	19	7.15	3.626	upstream	1931	PTGER1	

Output Folder Gene Name ER View in IGV close

Figure 7: Annotation table. This window exhibits annotated peaks. It contains peak location, nearest gene and distance to TSS. This window also provides, easy access to IGV for visualizing peaks and automatically generated normalized ATAC-seq signal by GUAVA. One can search peak by symbol of the nearest gene using search box provided at the bottom.

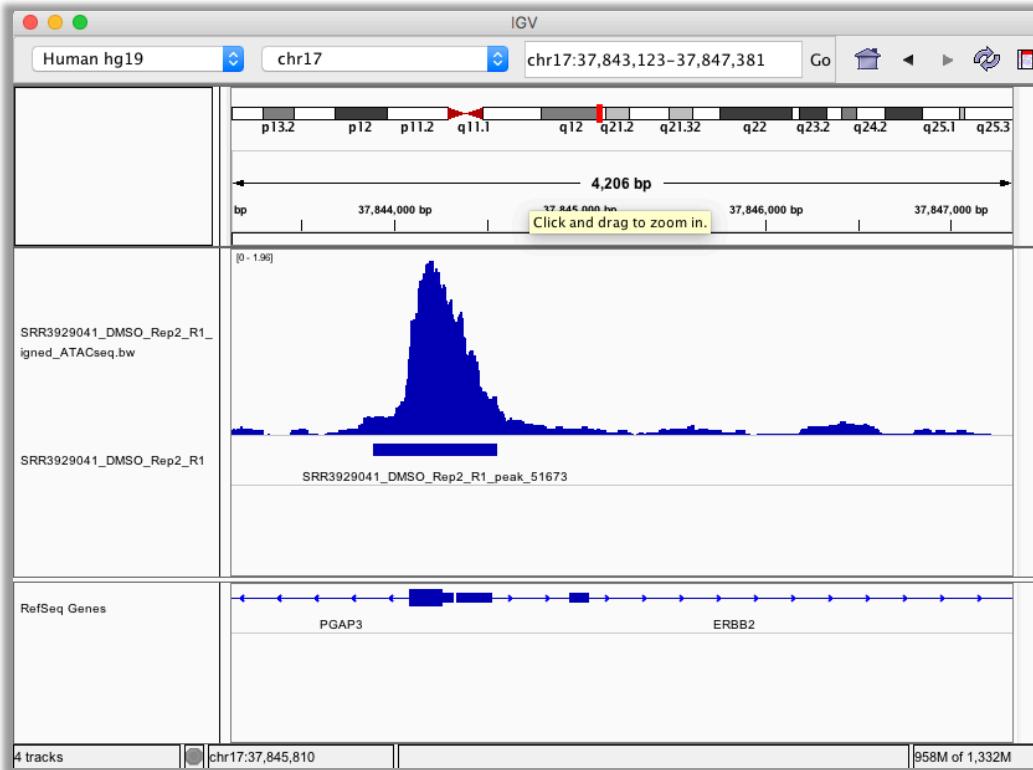


Figure 8: Visualization of ATAC-seq peaks with IGV.

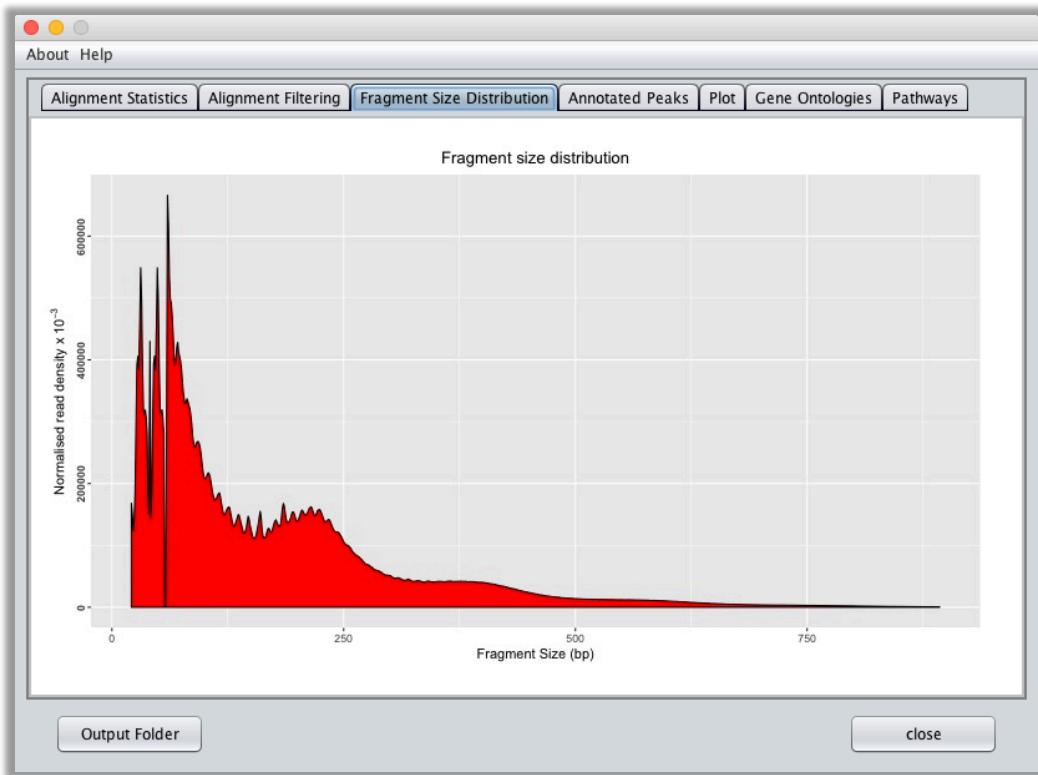


Figure 9: Graph showing the fragment size distribution.

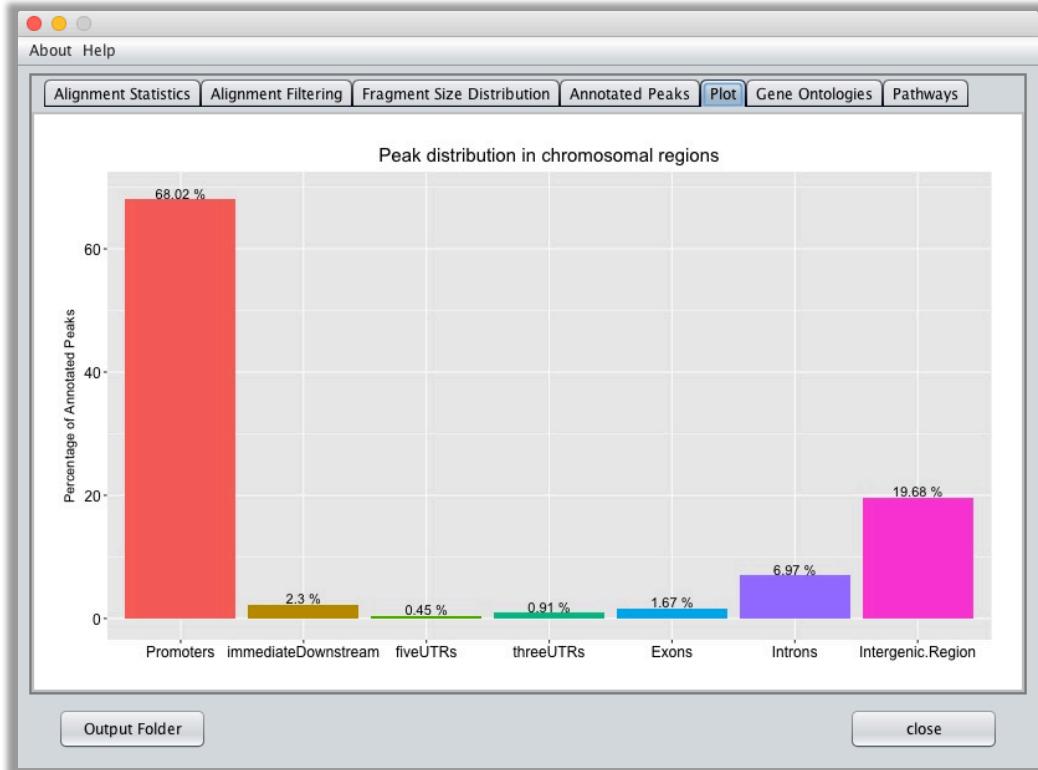


Figure 10: Bar chart showing the percentage of peaks in various genomic locations such as promoter, intron, exon, UTR, etc.

Figure 11 shows a screenshot of a software interface for analyzing gene ontology terms. The window title is 'Gene Ontologies' and it contains tabs for Alignment Statistics, Alignment Filtering, Fragment Size Distribution, Annotated Peaks, Plot, Gene Ontologies, and Pathways. The 'Gene Ontologies' tab is selected. A table lists over-represented GO terms with columns for GO ID, GO Term, Type, P value, adj. P value, and Gene Symbols. The table includes rows for various biological processes like 'pattern specification process', 'purine ribonucleoside triphosphate ...', and 'regulation of macromolecule biosynt...'. At the bottom left is an 'Output Folder' button, and at the bottom right is a 'close' button.

Figure 11: Over represented gene ontology terms.

Figure 12 shows a screenshot of a software interface for analyzing pathways. The window title is 'Pathways' and it contains tabs for Alignment Statistics, Alignment Filtering, Fragment Size Distribution, Annotated Peaks, Plot, Gene Ontologies, and Pathways. The 'Pathways' tab is selected. A table lists over-represented KEGG pathways with columns for KEGG ID, Pathway Name, P value, adj. P value, and Gene Symbols. The table includes rows for various biological pathways like 'ErbB signaling pathway', 'Spliceosome', and 'MAPK signaling pathway'. At the bottom left is an 'Output Folder' button, and at the bottom right is a 'close' button.

Figure 12: Over represented pathways.

Furthermore, these results are stored in the output folder, click the ‘output folder’ button at the bottom-right to open the output folder.

ATAC-seq differential analysis program

This program compares ATAC-seq signals from two conditions and returns the differentially enriched signals. Additionally, it provides the peak annotation and functional analysis for differentially enriched peaks. There are two input windows for this program. First window is to upload the ATAC-seq signals from different conditions and replicates (Figure 13 and 14). Second window allows you to specify differential analysis related parameters e.g. fold change (Figure 14).

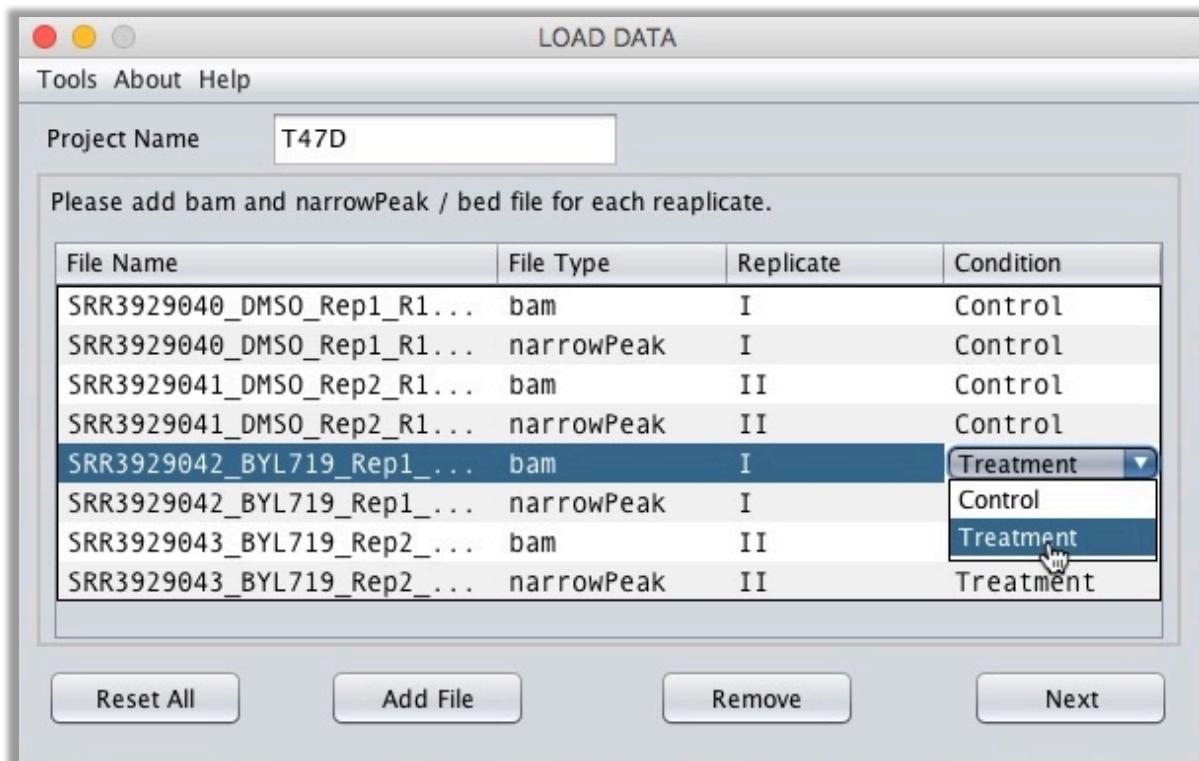


Figure 13: GUAVA ATAC-seq differential analysis input interface 1

Use 'add file' and 'remove' buttons to add and delete input files respectively. Once you have uploaded bed file containing ATAC-seq peaks and bam files, specify the condition and replicate number for each file. Click 'Next' to specify differential analysis parameters.

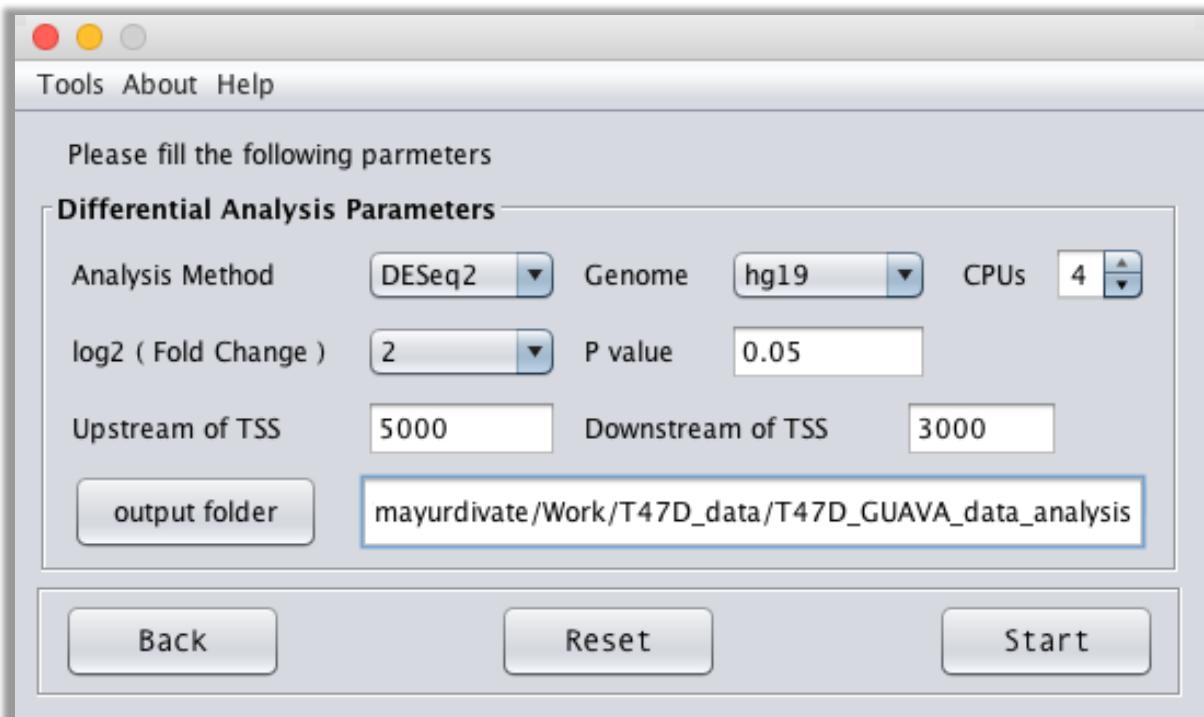


Figure 14: GUAVA ATAC-seq differential analysis input interface 2

Once you have filled all the required parameters click ‘Start’ button to run differential analysis.

Analysis method: select method for differential analysis DESeq2.

log2 (Fold Change): log2 fold change cut off to define differentially enriched peaks. Default 2.

P value: P value cut off to select most significant differentially enriched peaks. Default 0.05.

Upstream of TSS: if the peak is present within a specified distance (in base pair) from the TSS of a gene, to the upstream. Then that gene will be associated with the peak for functional analysis. Default 5000.

Downstream of TSS: If the peak is present within a specified distance (in base pair) from the TSS of a gene, to the downstream. Then that gene will be associated with the peak for functional analysis. Default 3000.

Output folder: Select folder to save GUAVA differential analysis results.

Output interface of ATAC-seq differential analysis program

The output interface of ‘ATAC-seq differential analysis’ program is also tabular like ‘ATAC-seq data analysis’ program.

- 1) ‘Summary’ tab: This tab provides summary of input parameters e.g. fold change cut-off, list of input files used for differential analysis (Figure 4A) etc.
- 2) ‘Differential Table’ tab: This provides the list of differentially enriched ATAC-seq signals with annotation such as nearest gene to peak and the distance between them (Figure 4B) etc. Same as output interface of ‘ATAC-seq data analysis’ program, there is a search box and ‘view in IGV’ button at bottom of window. Which can be used to sort peaks by gene symbol and view peaks in IGV from input samples, respectively (Figure 4D).
- 3) ‘Plot’ tab: This provides volcano plot of differentially enriched peaks (Figure 4C).
- 4) ‘Go Analysis’ and 5) ‘Pathway Analysis’ tabs: These tabs provide results of functional analysis i.e. enriched gene ontologies (Figure 4E) and pathways (Figure 4F) respectively.

Figure 15: Input summary.

Figure 16: Differentially enriched peaks with sorting and filtering functionality. Easy access to IGV to visualize differentially enriched peaks and normalized ATAC-seq signals from each sample.

Figure 17: Volcano plot indicating the differentially enriched peaks. Red: peaks with increased chromatin accessibility, green: peaks with reduced chromatin accessibility and black: peaks with no significant change in chromatin accessibility.

Figure 18: PCA plot

Figure 19: Peak visualization in IGV.

Figure 20: over-represented gene ontologies and

Figure 21: over-represented pathways.

How to download genome fasta file?

Fasta is a flat file format for representing nucleotide or protein sequences. Genome fasta file is a fasta file which contains the nucleotide sequences from all of the chromosomes of a particular organism. Genome fasta file is required for read mapping using any aligner tool. For users who don't know where they can use the UCSC link given below and choose desired organism to download fasta file,

UCSC link: <http://hgdownload.soe.ucsc.edu/downloads.html>

Then, click on the 'full data set'. This will open a new page, scroll down and click on the chromFa.tar.gz to download genome sequence.

Use following command to extract chromosome files and merge them into one file.

```
tar -zxvf -d /path/to/chromFa.tar.gz  
cat chromFa/*fa > GenomeBuild.fasta
```

How to create a genome index from genome fasta file

It is true that the genome fasta file is required for the alignment. But the aligners use special set of files called as genome index, generated from genome fasta. Index files are used to speed up the read mapping process so that the aligner can map millions of reads within few hours of time. Therefore, you need to create genome index file before read mapping. Remember that the genome index format is different for each aligner. Please refer to 'Download genome fasta file' section to find more information about downloading genome fasta file. If you already have a genome fasta file, then use Genome index builder program to create genome index. This program just takes genome fasta file and output folder as an input.

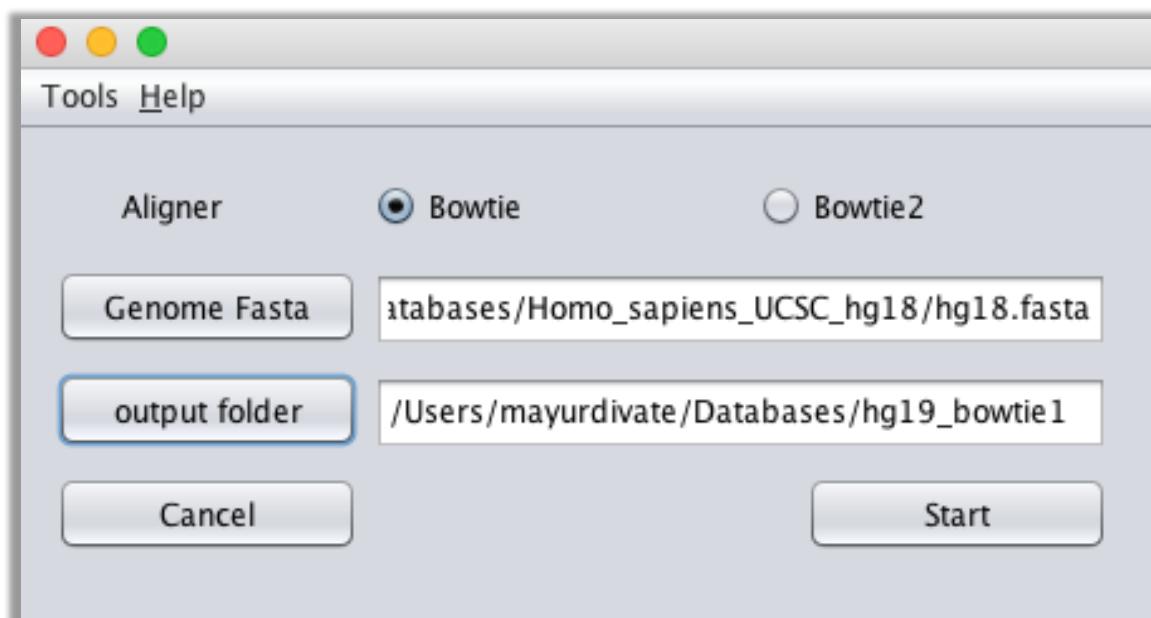


Figure 22: GUAVA Genome index builder program.