# Contents

目录

# Chapter 1

# 1 Introduction



## 1.1 What is HEMUdb?

After reviewing the existing literature on Andropogoneae, it was evident that a comprehensive database for the primary species within the Andropogoneae group was lacking. Consequently, we developed HEMUdb, a user-friendly and comprehensive database tailored for researchers

focusing on Andropogoneae.

## 1.2   What do we provide?

We have developed a user-friendly graphical interface that allows users to effortlessly explore genomic data across various representative Andropogoneae species with just a click. Leveraging a total of 4287 RNA-seq datasets from the public database, we applied a widely recognized transcriptome analysis process to construct 73 significant genomes of the Andropogoneae tribe.

We have designed four distinct search toolkits for users to utilize. Our server infrastructure is built on Django, MySQL, and Shiny technologies. You can conveniently access various types of visual analysis results (further details provided below). Should you have any inquiries or suggestions regarding our HEMUdb, feel free to contact us via email at **zhuyzh37@mail2.sysu.edu.cn** (Edward Zhu).

## 1.3   A useful function — Task ID

For enhanced user convenience, a distinctive Task ID will be assigned during the execution of the process. This Task ID will be visible in the URL and will remain effective indefinitely. Therefore, you have the option to save the Task ID and retrieve its associated results whenever needed.

A task has been submitted to the queue.
Unique task ID:

c291dbd7-0c1e-4875-8fc2-935ebe0a7824

**Permanent Task-ID**

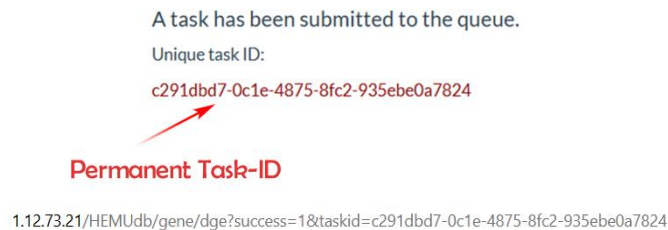1.12.73.21/HEMUdb/gene/dge?success=1&taskid=c291dbd7-0c1e-4875-8fc2-935ebe0a7824

Figure 1.3.1 an example of Task-ID and its place on URL

# Chapter 2

# 2   Toolkit I: Genome

## 2.1  Gene Information & Structure Search

We have developed a query page using the Shiny framework, allowing you to explore the gene structure of representative Andropogoneae species. To access this information, follow these straightforward steps:

- **Enter the gene ID you're interested in.**
- **Choose the relevant species.**
- **Click the Search button and wait for the results to load.**

This will provide you with details about the major transcript, CDS (Coding Sequences), UTR (Untranslated Regions), and other structural information associated with the specified gene.



Figure 2.1.1 interface of Gene Structure Visualization Search

You can track the query's progress and results loading by observing the progress bar located at the bottom right corner of the page. Once the loading is finished, the generated results will be displayed on the same page.
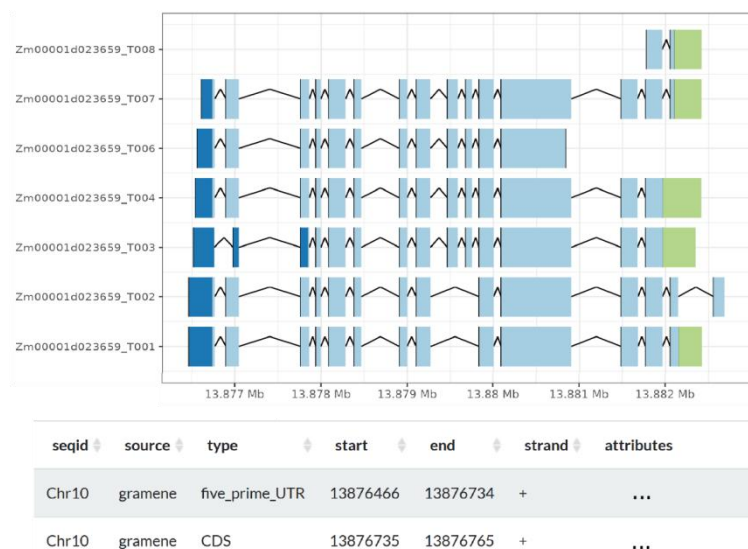


Figure 2.1.2 results of Gene Structure Visualization and details of each gene element

The initial figure in the results will provide a graphical representation of the primary transcript details for the gene, with distinct gene structures represented by various colors. The number of

displayed transcripts can be modified using the option located in the lower left corner. Subsequently, a table will present comprehensive information about each gene element. This information encompasses attributes such as source, start/end positions, and more.

## 2.2  Gene Functional Annotation Search

We have analyzed the gene functional annotation data for selected Andropogoneae species. As a result, we've created a user-friendly interface to search through the annotation outcomes. You can explore gene functional annotations across RNA-seq samples that are part of the HEMU catalog. To do this, you'll need to use the specific Gene ID for each genome. Here's how you can proceed:

- **Enter the Gene ID(s) you are interested in (separated by semicolons).**
- **Choose the relevant species from the provided options.**
- **Click on the "Search Gene Function" button and wait for the results to appear.**



Figure 2.2.1 interface of Functional Analysis

After a brief waiting period of a few seconds, the results will be presented within the same interface. A table will display comprehensive information, encompassing the primary transcript ID, a concise functional description, as well as GO/KEGG term and KEGG pathway term details.

### Gene function query: result

| gene ID | canonical transcript ID | function description | GO term | KEGG term | KEGG pathway term |
|---------|------------------------|---------------------|---------|-----------|-------------------|
| CI017000 | CI017000_T3 | callose synthase | GO:0005575,GO:0005623,GO:0005886,GO:0016020,GO:0044464,GO:0071944 | ko:K11000 | - |
| CI017097 | CI017097_T2 | Pectinesterase | GO:0005575,GO:0005618,GO:0005622,GO:0005623,GO:0005737,GO:0030312,GO:0044424,GO:0044464,GO:0071944 | ko:K01051 | ko00040,ko01100,map00040,map01100 |
| CI017258 | CI017258_T3 | Epsin N-terminal homology (ENTH) domain | GO:0005575,GO:0005623,GO:0005886,GO:0016020,GO:0044464,GO:0071944 | ko:K20043,ko:K20044 | - |

Figure 2.2.2 results of Gene Function Query

## 2.3  Multi-omics Genome Browser

Multi-omics analysis is a systems biology approach that comprehensively explores molecular-level information within organisms by integrating various high-throughput technologies. This holistic method enables researchers to gain a more comprehensive and

in-depth understanding of the interactions and regulatory relationships among various molecules within a biological system. This page is designed to showcase visualized genomic data for individual species that have been processed. A total of 9 species are available, and custom parameters can be applied to visualize Multi-omics genome data in a linear view. The genome assemblies of six species have reached the chromosome level. Click on the green link below each species name to view the corresponding visualization. For a more immersive experience, click the "Go fullscreen" button at the top of the page.



Figure 2.3.1 interface of Multi-omics Genome Browser

Taking the example of Chromosome 7 of *Zea mays*, genomic data, including gene annotation, transposable elements, and ATAC-seq, is presented on an interactive page. You can select specific genomic information of interest from the checkboxes on the right. The most recently chosen content will be dynamically generated at the bottom of the page. Clicking on genes within the page provides more detailed information. Additionally, you can use other action buttons at the top of the page to adjust the displayed windows.
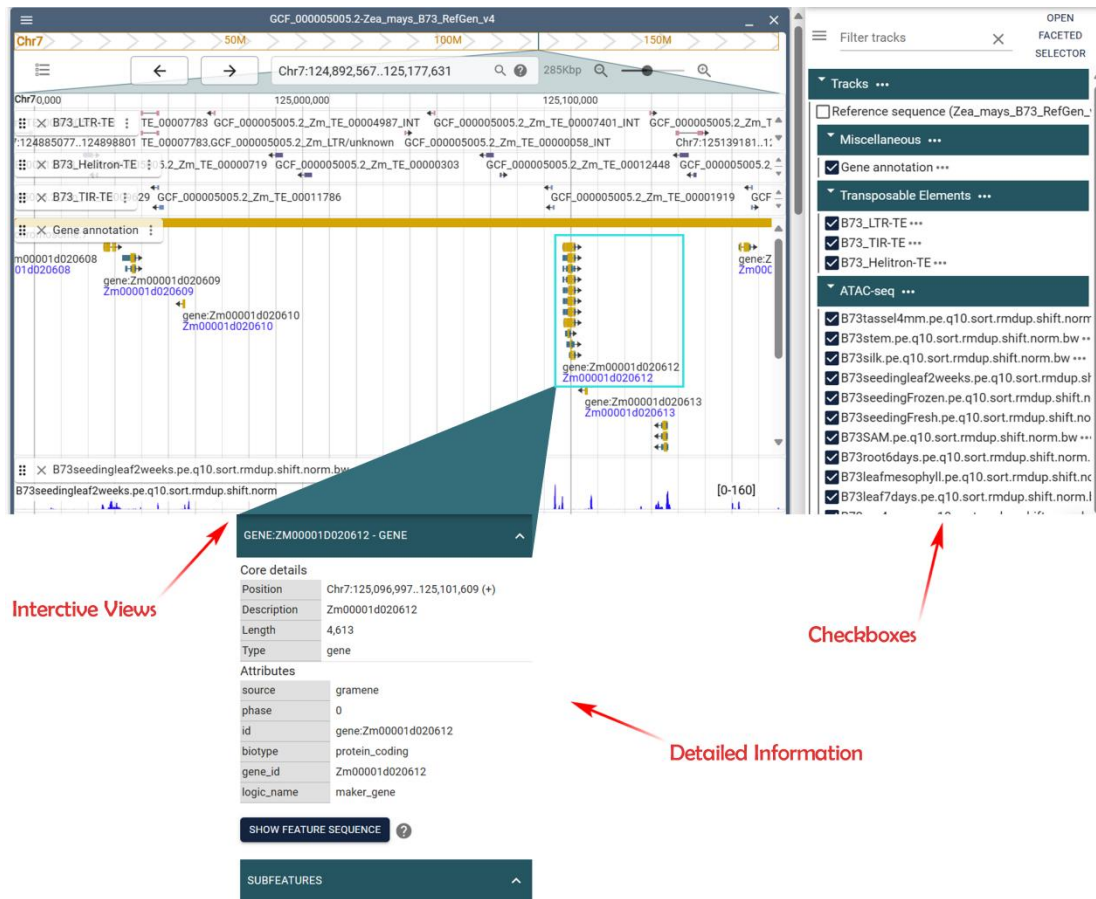
Figure 2.3.2 interactive interface of multi-omics genomic information

## 2.4  Genome Synteny Viewer

Genome synteny refers to the phenomenon of the conservation of the order and orientation of genes or other genetic elements across the genomes of different species. Synteny provides insights into the evolutionary relationships between species and can reveal the conservation of genetic content and organization. Currently, we have completed synteny analyses among three pairs of species and have visualized the results for user convenience. Click on the green link below each species to view the synteny results for that specific pair of species.
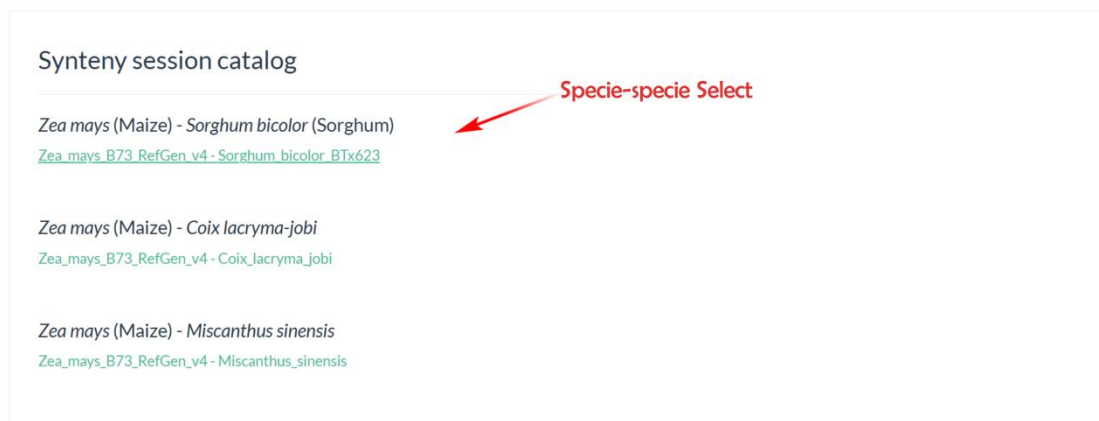
Figure 2.4.1 interface of Genome Synteny Viewer

Taking the example of chromosome 7 in *Zea mays* and chromosome 2 in *Sorghum bicolor*, in addition to various omics data such as genome annotation, transposable elements, ATAC-seq, etc., a synteny plot showing the corresponding positions of similar sequences between the chromosomes of the two species will be generated. Additionally, a set of Synteny anchors will also be generated by default. The framework used by viewers is the same as mentioned earlier, and relevant operational functions can be referred to in the methods described in the previous section on the **Multi-omics Genome Browser**.
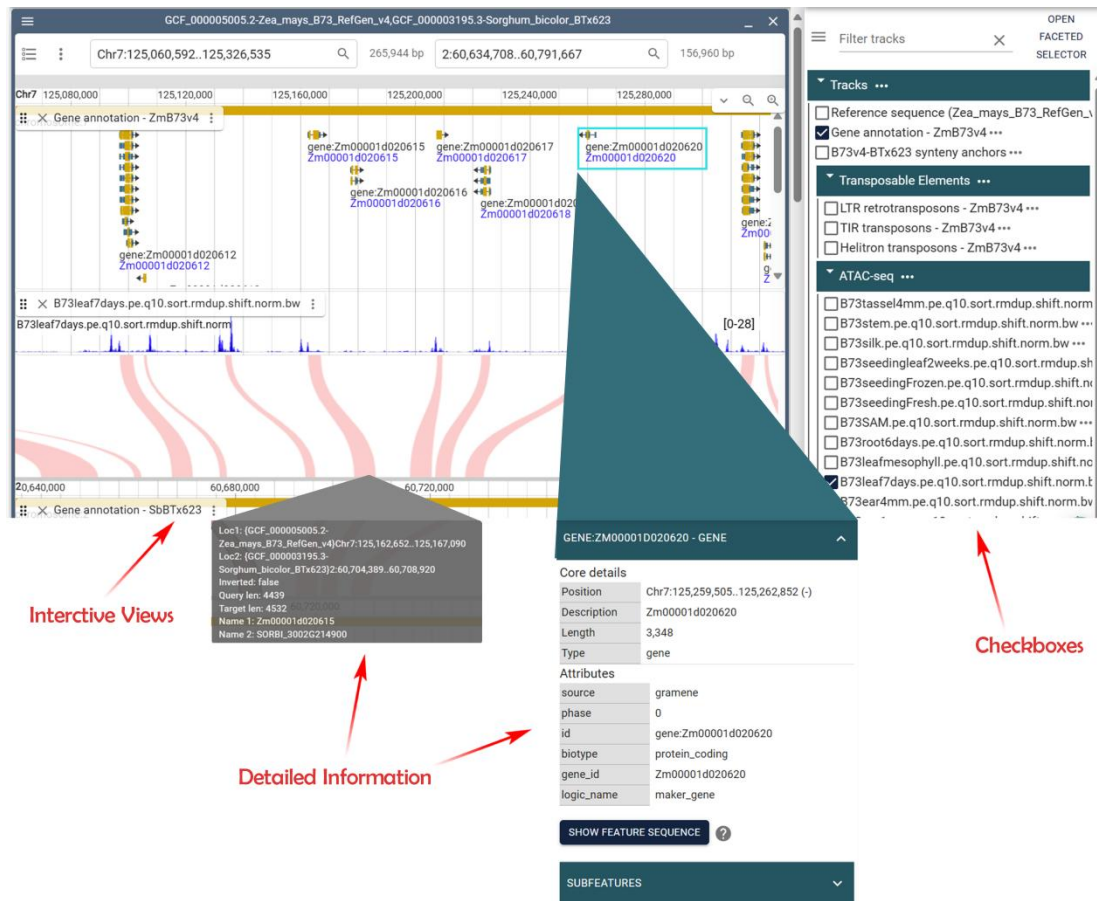


Figure 2.4.2 interactive interface of specie-specie genomic data information

## 2.5 The BLAST Server

We have incorporated part of the BLAST search engine framework from Sequance Server (sequenceserver.com) and supplemented it with nucleotide and protein databases assembled for our species. You can either directly access the page to input sequences of interest for analysis or, alternatively, click on "Send sequence to BLAST module" from almost any page that generates sequences to be redirected to this analysis page. By selecting the relevant database information, you can obtain alignment results.

Figure 2.5.1 interface of BLAST module

# Chapter 3

# 3  Toolkit II: Transcriptome

## 3.1  Gene Expression Profiles

We have introduced the **Online Query Module** for gene expression in the Andropogoneae transcriptome database. The search interface is shown in Figure 3.1.1. To obtain gene expression profiles, follow these steps:

 • **Enter the Gene IDs you are interested in into the left input box (use ";" to separate multiple samples).**
 • **Select the desired Andropogoneae family from the options on the right.**
 • **Choose your preferred expression format (FPKM and TPM are both supported).**

If you are using this function for the first time, you can click on the "Show Example" button. This will automatically fill in a prepared Gene ID in the left box.

Figure 3.1.1 interface of Gene Expression Profiles

After a brief wait of a few seconds, you will be able to access the results, presented through two panels showcasing gene expression profiles for each Gene ID.

**Expression Plots** We offer two interactive expression plots to visualize the profile outcomes, encompassing both sample-level and tissue-level data. In the sample-level plots, by hovering your cursor over any sample data point, you can view its corresponding ID, along with the associated TPM/FPKM value and sample ID. Likewise, within the tissue-level plots, you can employ the same approach to retrieve information about the tissue type and its corresponding TPM/FPKM value for any given sample data point.

**Fundamental Info** This section provides essential information about the queried gene, including frequency, maximum, minimum, and median expression levels across samples. This furnishes you with a fundamental understanding of the Gene ID you're interested in. By clicking the "Search gene sequence" button, you will be directed to the Raw Sequence Acquisition page, enabling you to access sequences for the target gene based on genome annotation profiles. Additionally, you have the option to preview or download the raw gene expression data table, which was used to generate the plots. This table supplies Sample IDs, TPM and FPKM values, as well as tissue types.

Figure 3.1.2 Expression Plots and Fundamental Info

## 3.2  Raw Sequence Acquisition

On this page, you have the option to swiftly retrieve gene, transcript, or protein sequences associated with your chosen Gene ID. Simply enter the Gene ID of interest into the query box. If you wish to query multiple Gene IDs simultaneously, please separate them using a semicolon ";".



Figure 3.2.1 original interface of Raw Sequence Acquisition

After a short wait of a few seconds, you will find the FASTA-formatted sequences displayed below on the page. There's no need to copy the sequences manually. Simply click the **"Send sequence to BLAST module"** button, and these sequences will automatically populate the BLAST window for you.
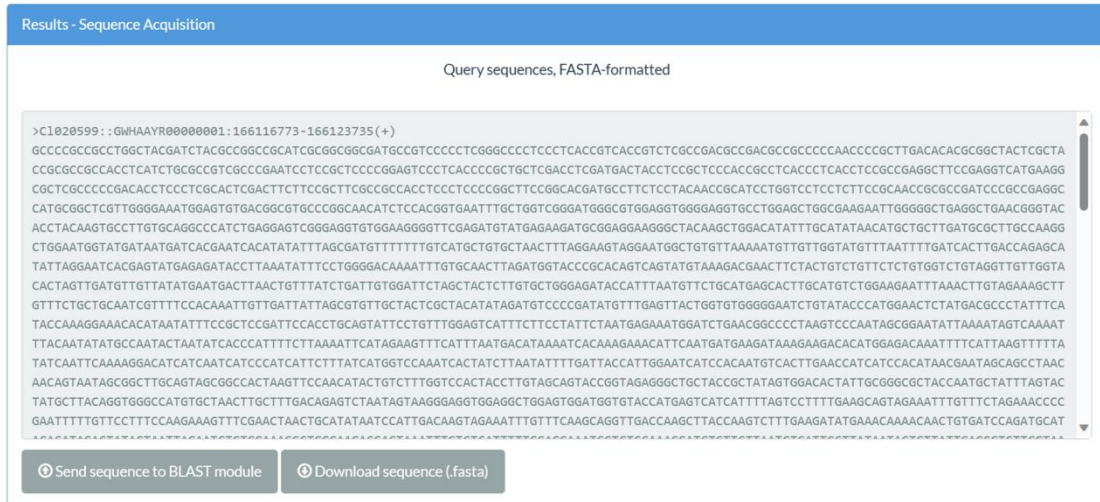
Figure 3.2.2 results interface of Raw Sequence Acquisition

## 3.3   Differential Gene Expression(DGE) Analysis

Once you have obtained gene expression data of interest, your next step might involve investigating the variations in expression across different tissues or organs. To assist with this, we offer a comprehensive platform for Differential Gene Expression (DGE) Analysis.
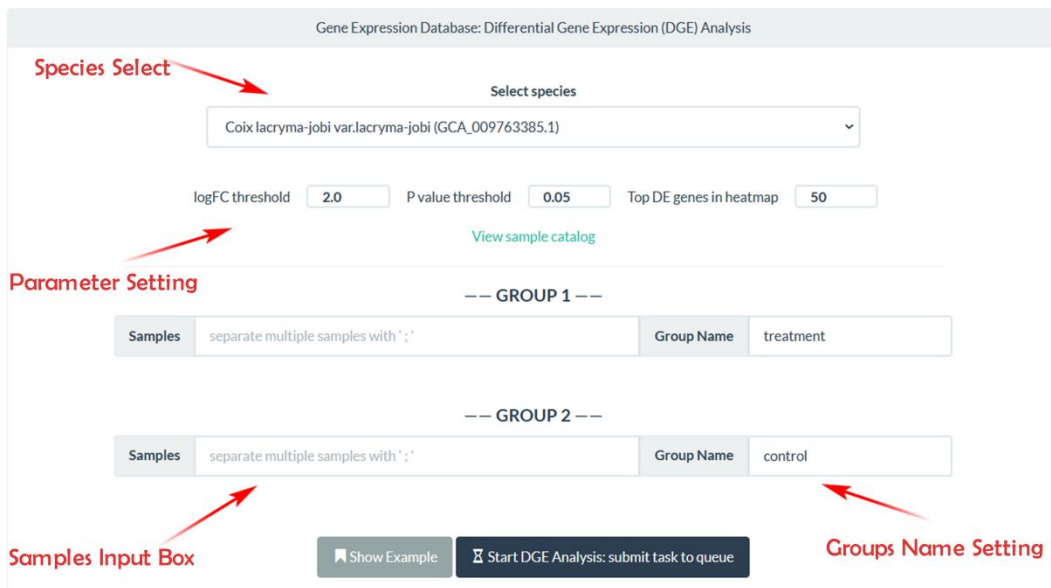


Figure 3.3.1 the interface of Differential Gene Expression Analysis

Here's a step-by-step guide to help you navigate the process:

· **Select the species to which your genes belong.**
· **Adjust the relevant parameters until they are appropriate (default settings are also available).**
· **Enter the series of samples into the provided input box, separating them with**

**semicolons.**
- **Assign brief names to the groups for easier identification in the subsequent results.**
- **Initiate the analysis by clicking the "Start" button.**

If you are using this function for the first time, you can click on the "Show Example" button. This will automatically populate the input box with a set of prepared sample sequences.

After waiting for approximately 2 minutes, a comprehensive set of analysis results will be presented. These results have been categorized into four distinct sections for ease of comprehension.

Project Summary    n this section, we provide an overview of the species, groups, and sample accessions you have selected for analysis. Additionally, we have compiled a table in CSV format that highlights the differentially expressed genes between groups, accompanied by detailed parameters.

Overview & Normalization of Expression Data    This section includes three figures that elucidate the expression data provided. The first figure features raw and normalized sample TPM (Transcripts Per Million) boxplots, showcasing the distribution before and after normalization for each sample. The second figure depicts density curves of gene expression across samples, offering insight into the overall expression trends of genes. The third figure presents a two-dimensional correlation heatmap of expression among samples, using varying shades of color to visually represent the correlation coefficients between samples.

Principal Component Analysis    In this section, we offer a Principal Component Analysis (PCA) plot. This plot aids in intuitively grasping the similarity between samples. For result accuracy, ellipses are excluded when the number of samples is insufficient.

Differential Analysis    his part features a classical volcano plot, which effectively displays the differential expression status between groups. By referring to the provided key, you can readily understand the distribution of genes categorized as up-regulated, down-regulated, or not significantly altered. Subsequently, a Heatmap showcasing the most differentially expressed genes (based on TPM) illustrates the overall expression changes across multiple samples.

Preview and download options are made available. You can simply click the corresponding button beneath each project to preview or download the results as needed.

## 3.4  GO/KEGG Enrichment

We offer a user-friendly tool for conducting GO/KEGG enrichment analysis. This tool enables you to explore functional annotation details of gene transcripts in representative Andropogoneae species. To proceed, simply follow these steps:
- **Gather the gene IDs you're interested in.**
- **Enter the gene IDs in the designated Gene List box (each ID separated by a line break).**

- **Select the enrichment method you need (GO/KEGG).**
- **Click the Enrich button to submit your request.**



Figure 3.4.1 interface of GO/KEGG Enrichment

If you are using this function for the first time, you can click the **"Show Example"** button. A selection of gene lists that we have prepared will be automatically displayed in the left box. (Currently, we have included effective data and an example from *Zea mays*.)

Once the analysis is complete, a figure and a table will be presented. The figure is an interactive bubble plot where the color of each bubble corresponds to its p-value, and the size of the bubble represents the number of genes it contains. The ontology associated with each bubble is displayed on the right side of the axis. Hover your cursor over any bubble to view detailed information about it. Additionally, you can click the button below to download a static image of the plot.



Figure 3.4.2 bubble plot of GO/KEGG Enrichment result

The table contains comprehensive data from the Enrichment Analysis, and it includes options for previewing and downloading the data.

## 3.5  Weighted Gene Co-expression Network Analysis(WGCNA)

Weighted Gene Co-expression Network Analysis (WGCNA) is a commonly used analytical method in systems biology and bioinformatics. This method is based on the expression patterns of genes across different samples, constructing a co-expression network among genes. By calculating the correlation between pairs of genes, highly correlated genes are grouped into modules. Each module represents a set of genes with similar expression patterns across samples. This approach not only focuses on the expression changes of individual genes but also emphasizes the coordinated regulatory relationships among sets of genes. It helps researchers gain a better understanding of the relationship between genes and phenotypes.

We have divided the overall WGCNA process into three sections. Section 1 covers Data curation, filtering, and sft selection. Section 2 involves Co-expression network construction, and Section 3 focuses on Module-trait correlation analysis. It's noteworthy that the analyses in Section 2 and Section 3 are relatively independent, as you will observe from the subsequent workflow. The detailed operational steps are as follows:

 • **Navigate to the Section 1 interface, select the species you wish to analyze, and input the Sample Accession List and Gene List into the corresponding fields after adjusting relevant parameters**

 • **Click the "Submit Task" button to initiate the analysis**

 • **After the initial filtering, you will receive a unique Project ID. This ID will allow you to view completed analysis results at any time from the Result Viewer. The Section 1 results page includes a Program log file, where the "SFT prediction" will assist you in Section 2 analysis.**

 • **Switch to the Section 2 page, input the Project ID and Recommended SFT soft power into the corresponding fields, and click the button to perform the analysis**

 • **Switch to the Section 3 page, input the Project ID into the corresponding field, and click the button to initiate the analysis**

Figure 3.5.1 interface and flow path of WGCNA

After all analyses are completed, you can open the Result Viewer interface at any time. Enter the Project ID obtained during the analysis to quickly view or download the analysis results. The entire process will generate four figures: Sample cluster dendrogram, Scale independence and mean connectivity plot (from section 1), Network and cluster dendrogram (from section 2), and Module-trait correlation heatmap (from section 3), along with the associated tables.



Figure 3.5.2 interface of Result Viewer

# Chapter 4

# 4   Toolkit III: Gene Family

## 4.1  Family Identification(HMM & BLASTP)

In most cases, genes belonging to the same gene family exhibit noticeable similarities in both structure and function. To cater to your requirements for identifying shared domains and distinctions within the gene sequences of your interest, we offer two distinct gene sequence alignment methods. Below, you will find a concise overview of these methods.

**Hidden Markov Model(HMM)**   A statistical modeling method used for predicting gene sequences based on known ones. Utilizing training data from established gene sequences, it identifies common structural and feature traits within gene families. This allows the creation of gene models to predict unknown gene sequences.

**BLASTP**   A protein sequence alignment method that identifies similarities between unknown protein sequences and a pre-prepared protein database. It assigns scores to assess the level of similarity, indicating the likelihood of a new gene when a higher score is achieved.

Identifying gene families requires careful consideration of multiple factors, as relying on single methods can lead to errors. To enhance the accuracy and reliability of gene identification, cross-validation is recommended. Follow these steps for improved results:

- **Select the desired protein sequences (separated by species).**
- **From the HMMs menu, select the relevant gene model.**
- **Adjust parameters as needed (defaults are provided).**
- **Click the "Search Family Members" button.**



Figure 4.1.1 interface of HMM-based Family Identification

If you are using this function for the first time, you can click on the **"Show Examples"** button. This will automatically populate the input box with a preset model that we have provided.

The HMM-based results include a score-sequence table, which is organized from highest to lowest score, indicating the similarity of each potential gene with the specific model. Additionally, comprehensive domain annotations for each sequence and alignments for individual domains are presented. The identification process concludes with a statistical summary.
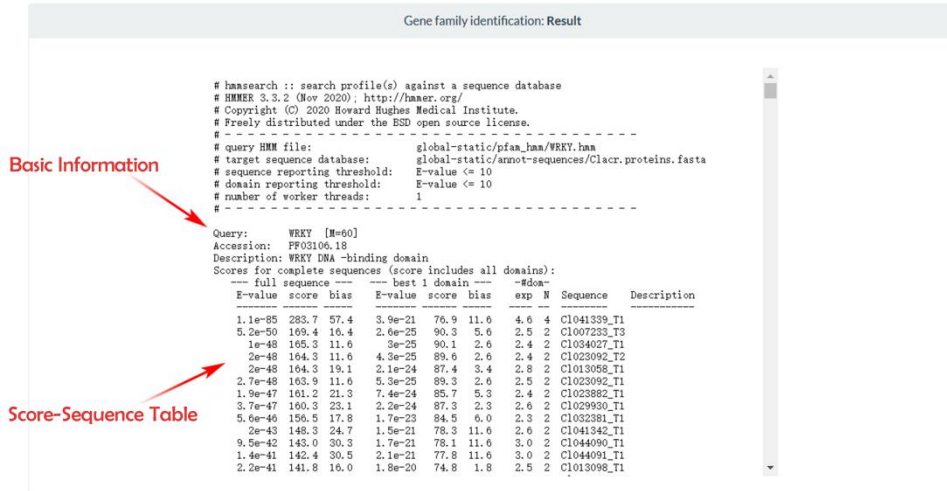


Figure 4.1.2 A example of result based on HMM

Regarding the outcome derived from BLASTP analysis, it is recommended to acquire a representative sequence (RefSeq) from a protein family to identify members in different species. The RefSeq can be obtained through manual curation or downloaded from a publicly accessible RefSeq database.
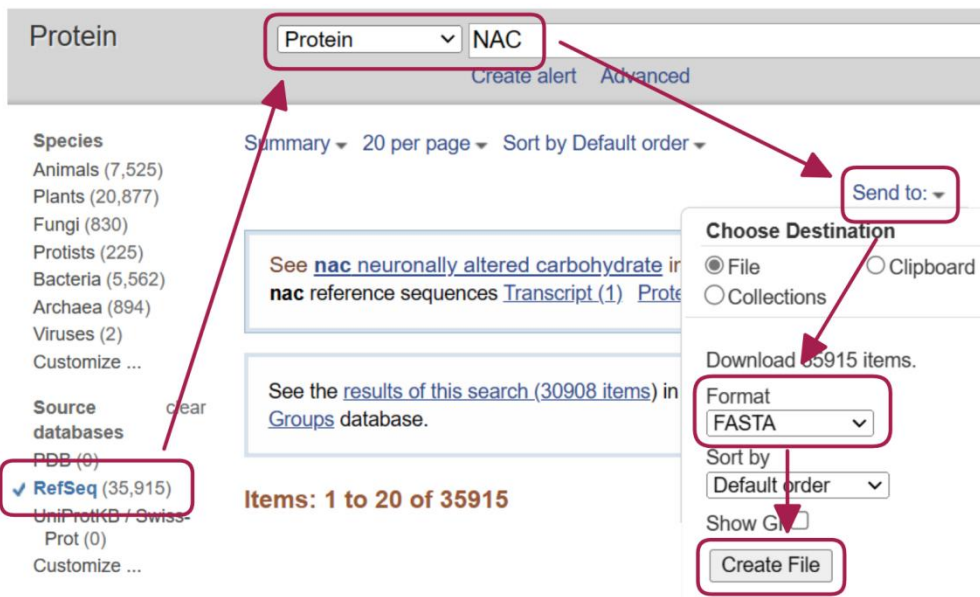
Figure 4.1.3 guide for downloading RefSeq from the NCBI database

To retrieve the relevant outcomes, kindly insert the preferred protein sequences within the designated input field. Following this, select the appropriate RefSeq, make adjustments to the advanced parameters (set to default if preferred), and proceed by clicking on the **BLAST** button. Subsequently, the comprehensive alignment details for the sequences will be presented in a systematic manner.

## 4.2 Phylogenetic Analysis

Once you have identified a set of potential genes through cross-validation, creating a gene family phylogenetic tree can help elucidate the evolutionary and genetic connections among these genes. To facilitate this, we offer a phylogenetic analysis toolkit for the automatic construction of a foundational tree. The recommended procedure involves the following steps:

- Select the desired species.
- Configure appropriate parameters and bootstrap replicates (default values are provided).
- Input the genes using either FASTA format or Gene IDs, obtained from your prior analysis.
- Click the **"Analyze"** button to submit the task to the queue.



Figure 4.2.1 the interface of Phylogenic Analysis

If you are using this function for the first time, you can click the **"Show Examples"** button. This will automatically populate the box with a set of preselected Gene IDs that we have prepared.

The outcome will present a preview of a phylogenetic tree, indicating the converted bootstrap value on each node. Generally, values above 50 indicate the reliability of this phylogenetic tree.

You have the option to download the static image in PNG format or obtain the editable tree in NWK format for importing into the MEGA (Molecular Evolutionary Genetics Analysis) software for further enhancement.

## 4.3 Family Expression Heatmap

For analyzing the expression levels of multiple gene families across various samples, we offer a heatmap generation tool designed to help you visualize your gene family data. Currently, this tool can handle a maximum of 100 sample accessions and 100 gene IDs for analysis. To utilize this tool, please follow these steps:

- **Select the species of your samples.**
- **Choose the preferred expression level format (defaulted to TPM).**
- **List your samples and gene IDs (up to 100) and input them in the provided boxes.**
- **Click the "Submit Task" button to initiate data submission.**

For first-time users, you can click the "Show Examples" button. This action will automatically populate the boxes with sample accessions and gene IDs from a set of prepared lists. (Currently, we have included effective data and an example from *Coix lacryma-jobi var. lacryma-jobi*.)



Figure 4.3.1 interface of Gene Family Expression Heatmap Generator

Consequently, the system will display a visual representation comprising a heatmap depicting gene expression levels and a dendrogram illustrating sample clustering. This will be accompanied by a comprehensive table showcasing the gene expression data across various samples.

Figure 4.3.2 visualization of gene family expression

Within the interface, you'll notice that the intensity of colors on the right side signifies the relative expression level of each gene family. A darker shade indicates a higher expression for that gene family. The precise expression is quantified by the TPM value. On the left side, a cluster dendrogram has been constructed based on the gene family's expression pattern. This dendrogram visually represents the similarities and differences between gene families. For a more in-depth analysis, you can download the table we offer, which contains comprehensive data for your examination.

# Chapter 5

# 5   Toolkit IV: TE Analysis

## 5.1   TE(Transposable Elements) Expression Query

TE expression data of representative Andropogoneae species has been analyzed and an interactive query interface has been provided. During our analysis, individual TEs are classified into families based on the 80-80-80 rule proposed by Wicker et al. That is, two elements belong to the same family if they share 80% (or more) sequence identity in at least 80% of their coding or internal domain, or within their terminal repeat regions, or in both. You can query any TEs

you interested in this interface. Concrete steps we suggest are following:

- Select a suitable query mode
- Input TE ID/Sample Accession ID you interested in the box
- Choose corresponding species and the format you want
- Click **Search Expression Profile** button to query



Figure 5.1.1 interface of TE Expression Query

Waiting for a few seconds, you can see the results with two panels of the gene expression profiles for every TE ID.

**Expression Plots** We provide two interactive expression plots to display the profile results, both sample level and tissue level. In Sample level plots, stay your cursor on any sample data, then you can see its ID, value of TPM/FPKM and sample ID. Similarly, in tissue level, you can query the tissue type and its value of TPM/FPKM any sample data belongs to by the same way.

**Fundamental Info** This panel displays fundamental information regarding the query gene. Including frequency, max, min and median expression in samples, which give a basic recognition of your interested TE ID. The standard of an expressed TE family is set to be TPM/FPKM>1. Additional, buttons of previewing or downloading raw TE expression data table used to generate the plots are provided.
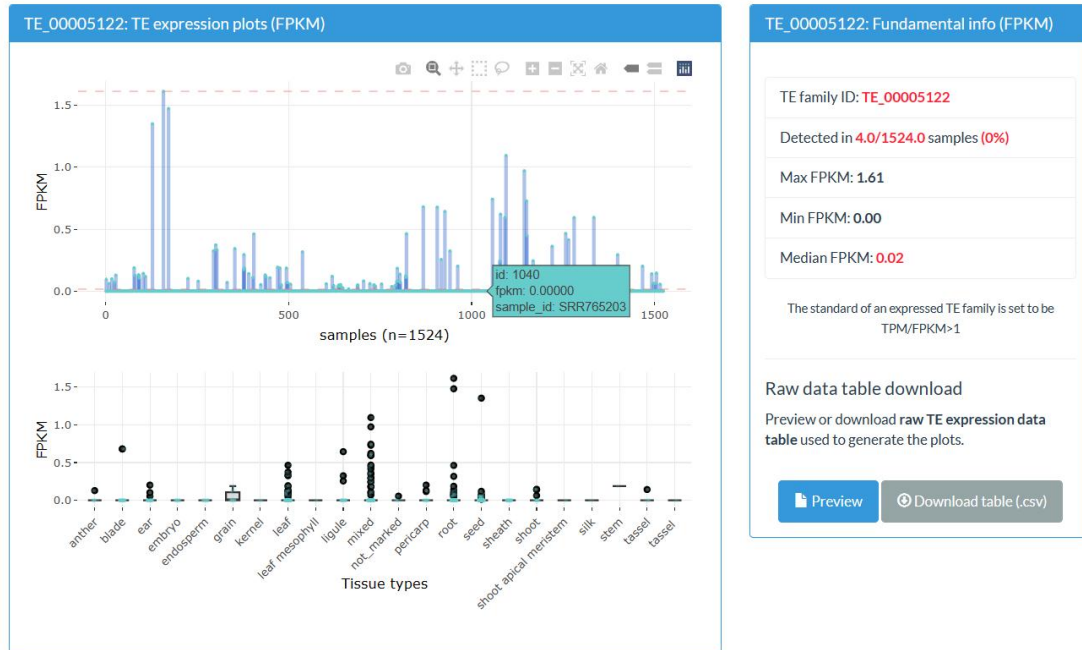
Figure 5.1.2 results of TE Expression Query

## 5.2 TE Insertion Location Search

TE(Transposable Element) is a type of DNA sequence that can move to different locations in the genome, and its insertion location in the genome is closely related to its biological function. To better understand the impact of TEs on genome structure and function, we have provided two methods for querying TE insertion locations: inputting Gene ID or sequence region range. These methods can help researchers identify where TEs are inserted, and further investigate their roles and biological functions in the genome. Concrete steps we suggest are following:

- Select the species you want to query
- Select suitable flanking region length(1.0 Kbp defaulted)
- Input the Gene ID or sequence region range you interested in the corresponding box
- Click **Search Relative TEs** button and wait for a minute

If you use this function first, you can click the **Show Examples** button. A preseted Gene ID/sequence region range will be filled in the corresponding box automatically.

Figure 5.2.1 interface of TE Insertion Search

About the results, a table of the information of insertion location in what you query will be provided, including its classification, start and end, score, p-value and so on.



Figure 5.2.2 results of TE Insertion Search(query by Gene ID)

Additionally, if you query it by inputting Gene ID, its detailed information will be enclosed.

## 5.3 Chromosome-level TE Insertion Density Search

An interactive Chromosome level TE Insertion Density Search platform is provided, which calculates the number of TE on chromosomes to determine their distribution and study their role in genome structure and function. At the top, species can be selected. Select the desired chromosome on the left and adjust the relevant parameters on the right to obtain the corresponding TE annotation data and draw example plot. Sliding the left slider can adjust the visible window size. Concrete steps we suggest are following:

- Select the species you want to query
- Set suitable parameters and sequences you interested
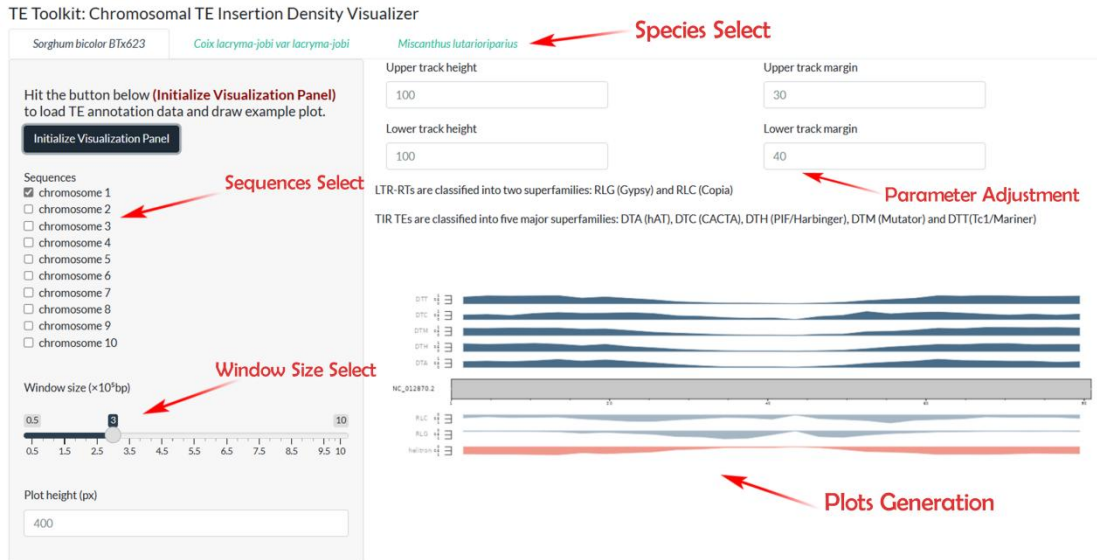- Click **Initialize Visualization Panel** button and wait for a minute



Figure 5.3.1 interface of TE Insertion Density Visualizer(interactive)

About the results, a series of TE insertion density example plots will be generated. Defaulted window size is $3 \times 10^5$ bp.

# Chapter 6

# 6   Toolkit V: Epigenome Analysis

## 6.1  Chromatin Accessibility Search

Chromatin Accessibility Search is an analysis method that identifies open chromatin regions based on peak information from annotated ATAC-seq samples. By identifying open chromatin regions in the genome, this method helps to understand gene expression regulation mechanisms and gene regulatory events during cell differentiation and development. Two ways to search for open chromatin regions are provided: inputting a Gene ID from a certain species or a sequence range on a chromosome, both of which will return relevant information on all open chromatin regions at that location. Concrete steps we suggest are following:

- Select the species you want to query

- Select suitable flanking region length(1.0 Kbp defaulted)
- Input the Gene ID or sequence region range you interested in the corresponding box
- Click **Search Relative Peaks** button and wait for a minute

If you use this function first, you can click the **Show Examples** button. A preseted Gene ID/sequence region range will be filled in the corresponding box automatically.



Figure 6.1.1 interface of Chromatin Accessibility Search

About the results, a table of the information of all chromatin regions in what you query will be provided, including its start and end, score, p-value and so on.



Figure 6.1.2 results of Chromatin Accessibility Search(query by Gene ID)

Additionally, if you query it by inputting Gene ID, its detailed information will be enclosed.

## 6.2  ChIP-seq Peak Annotation Analysis

We built a database for ChIP-Seq data of representative Andropogoneae species and streamlined

the entire process of ChIP-Seq peak annotation analysis. Specify ChIP-Seq sample accession ID then you can get corresponding analysis results. Pull down the menu selection at the bottom of this interface. Then you can browse all available ChIP-Seq samples. It lists more detailed information of every sample so that can assist you in choosing the sample ID you interested. Concrete steps we suggest are following:

· Pull down the ChIP-Seq samples menu

· Choose a sample ID you interested then input it in the box

· Select the corresponding species

· Set TSS(Transcription Start Sequence) region length the results finally display(in unit of bp, default provided)

· Set additional parameters(not recommended)

· Click **Start Analysis** button and wait for a minute

If you use this function first, you can click the **Show Examples** button. A preseted sample accession ID will be filled in the boxes automatically.



Figure 6.3.1 interface of ChIP-Seq Peak Annotation

The results of the analysis include five subgraphs:

**Heatmap**   This subgraph displays the distance of every ChIP-Seq near the TSS region, every colored short line corresponds to a sequence. Typically, there will be a sequence-enriched region near the TSS region.

**Density Map**   This subgraph visualizes the enrichment status of the ChIP-Seq data in the previous subgraph. You can visually observe a peak appearing near the TSS region. The shape of the peak may reflect the type of sequences. Typically, sharp peaks correspond to TFs(Transcription Factors), while broad peaks correspond to DNA methylation.

**Venn Plot**   Due to the distribution of known TSS regions, a single peak may have multiple potential attributes. For example, when two TSS regions are in close proximity (2~3kb), a peak between them may have both Intergenic and promoter attributes. Therefore, peaks with multiple similar attributes were clustered and a Venn plot was generated to describe the features of all ChIP-Seq data.

**Pie Plot**   This subgraph classified all peaks based on their locations, resulting in a clear pie chart reflecting the proportion of each peak type. You may want to pay particular attention to the proportion of promoter, which may reflect the proportion of TF sequences in all ChIP-Seq.

**TF Binding Loci Distribution Map**     This subgraph depicts the distribution (predicted results) of all possible TF binding sites near the TSS region, and distinguishes the relative distance between each binding site and the TSS region by color.
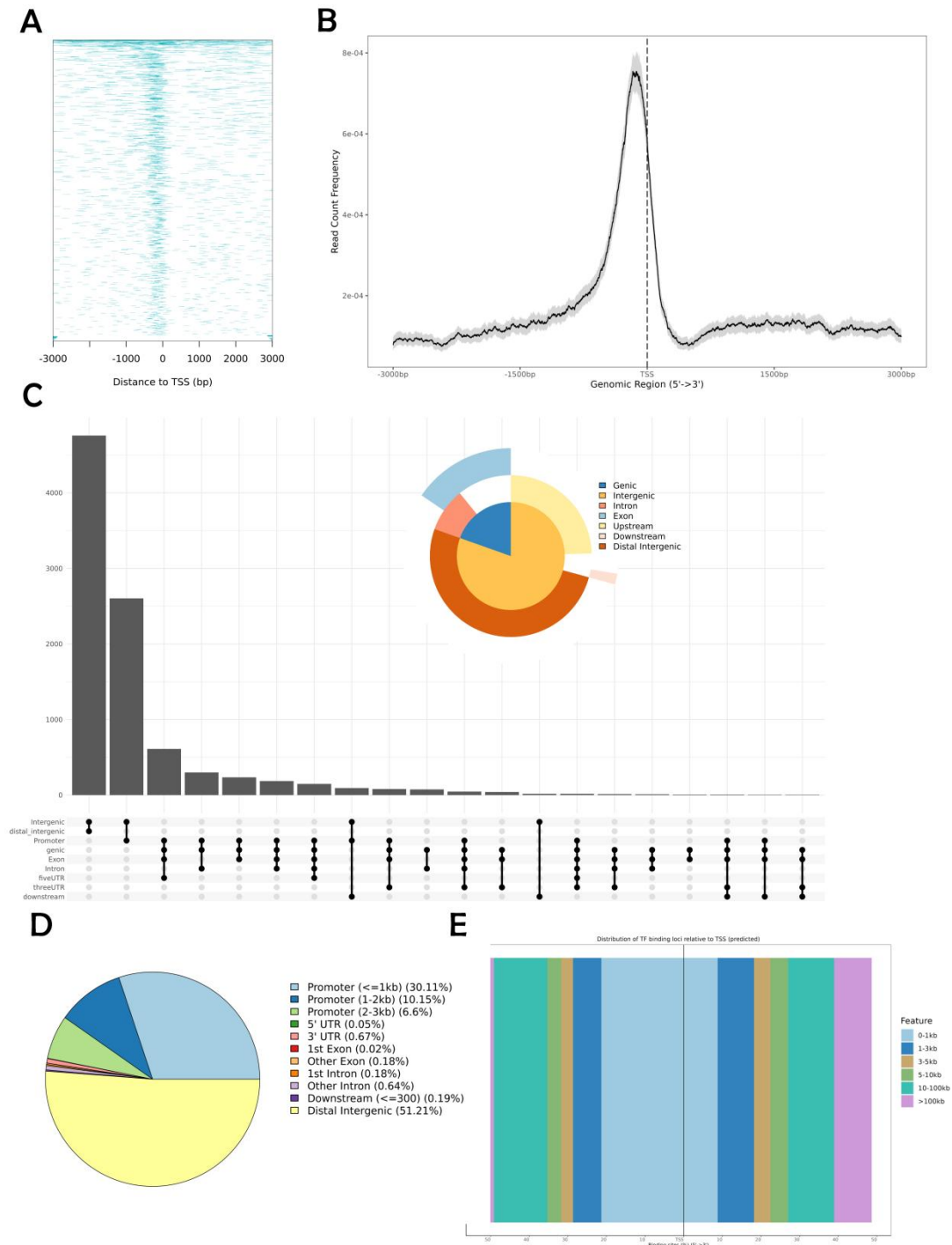


Figure 6.3.2 results of ChIP-Seq Annotation Analysis A. peak distribution heatmap B. peak distribution density map(interactive) C. proportion of peaks within certain genomic regions(venn plot) D. proportion of peaks within certain genomic regions(pie plot) E. distribution map of TF binding loci relative to TSS(interactive)

# Chapter 7

# 7    Data Warehouse

## 7.1  The HEMU user guide

To provide users with a concise understanding of the various **Toolkits** offered on this website, we have prepared a user guide. This guide outlines the overall framework of the website and provides instructions for using each module. You can view it online or download the PDF. As the website's features evolve, the content of the user guide will be adjusted accordingly.

## 7.2  Resources

All data used for display and visualization has been categorized into three sections: Genomic Resources, Transcriptomic Resources, and Epigenomic Resources. You can refer to these sections at any time to explore the data that interests you.