

mlPEA User Manual

(version 1.0)

- mlPEA is a user-friendly and multi-functionality platform specifically tailored to the needs of streamlined processing of m⁶A-Seq data in a reference genome-free manner. By taking advantage of machine learning (ML) algorithms, mlPEA enhanced the m⁶A-Seq data analysis by constructing robust computational models for identifying high-quality transcripts and high-confidence m⁶A-modified regions.
- mlPEA comprises four functional modules: **Data Preprocessing, Transcriptome Construction, m⁶A Calling, and Functional Exploration.**
- mlPEA was powered with an advanced packaging technology, which enables compatibility and portability.
- mlPEA project is hosted on <http://github.com/cma2015/mlPEA>
- mlPEA docker image is available at <http://hub.docker.com/r/malab/mlpea>

Data Preprocessing Module

This module provides four functions (see following table for details) to prepare epitranscriptome data.

Tools	Description	Input	Output	Time (test data)	Reference
Download File	Directly fetch epitranscriptome sequencing reads from NCBI's SRA database or other databases	SRR accession or HTTP/FTP link	Sequencing reads in SRA format	Depends on the network speed	SRA Toolkit
Sequence Data Preprocessing	Convert epitranscriptome sequencing reads from SRA to FASTQ format	Epitranscriptome sequencing reads in SRA format	Epitranscriptome sequencing reads in FASTQ format	~2 mins	SRA Toolkit
Assess Reads Quality	Check m ⁶ A-Seq reads quality and obtain high-quality reads	m ⁶ A-Seq reads in FASTQ format and adapter sequences in FASTA format	m ⁶ A-Seq reads in FASTQ format; MultiQC report in HTML	~2 mins	fastp ; MultiQC
Differential K-mer Analyses	Perform m ⁶ A-Seq differential k-mer analyses	m ⁶ A-Seq reads in FASTQ format	m ⁶ A-Seq reads in FASTA format	Depends on the file size	kmdiff

Download File

This function is designed to download epitranscriptome sequencing reads from NCBI SRA (Short Read Archive) database or from an user-specified HTTP/FTP link automatically. For the former, the **prefetch** function implemented in [SRA Toolkit](#) is wrapped to enable users to download sequencing data from NCBI SRA database; For the latter, **wget** command line is used to download the file according to an user-specified HTTP/FTP link.

Input

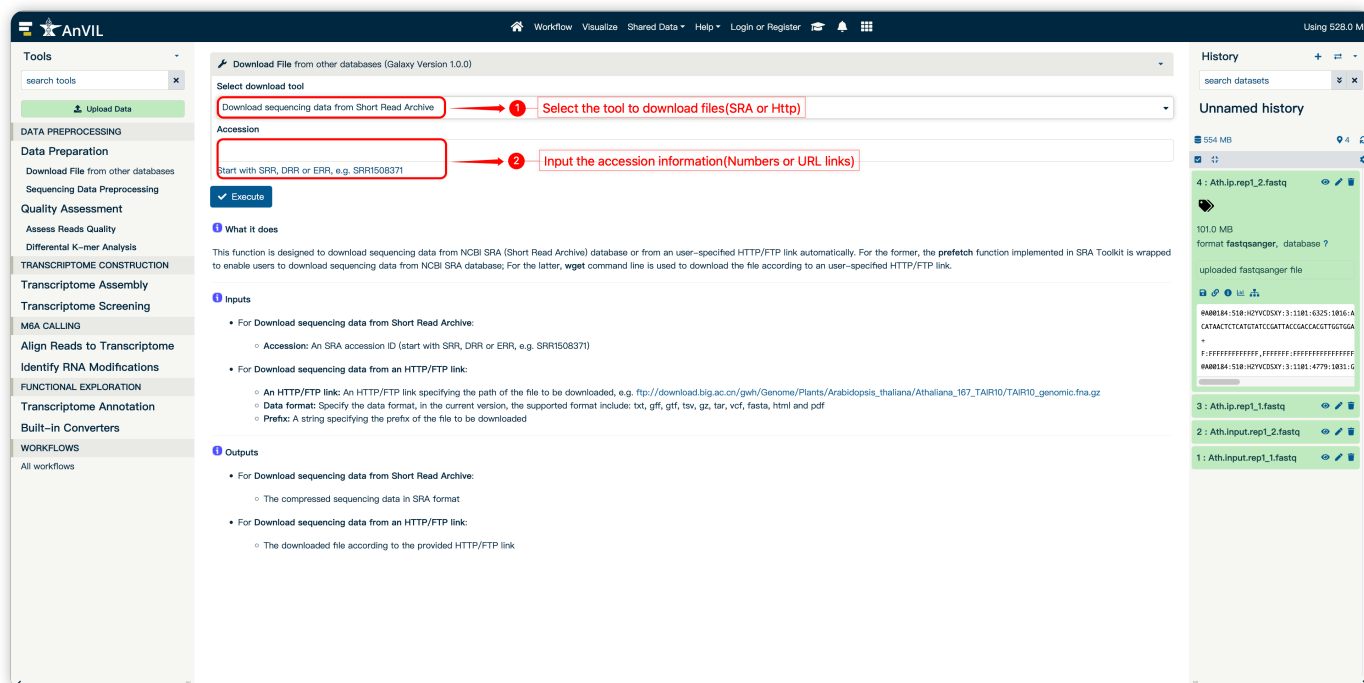
- For **Download sequencing data from Short Read Archive:**
 - **Accession:** An SRA accession ID (start with SRR, DRR or ERR, e.g. SRR1508371)
- For **Download sequencing data from an HTTP/FTP link:**
 - **An HTTP/FTP link:** An HTTP/FTP link specifying the path of the file to be downloaded, e.g. `ftp://download.big.ac.cn/gwh/Genome/Plants/Arabidopsis_thaliana/Athaliana_167_TAIR10/TAIR10_genomic.fna.gz`
 - **Data format:** Specify the data format, in the current version, the supported format include: txt, gff, gtf, tsv, gz, tar, vcf, fasta, html and pdf
 - **Prefix:** A string specifying the prefix of the file to be downloaded

Output

- For **Download sequencing data from Short Read Archive:**
 - The compressed sequencing data in SRA format
- For **Download sequencing data from an HTTP/FTP link:**
 - The downloaded file according to the provided HTTP/FTP link

How to use this function

- The following screenshot shows us how to download sequencing reads in SRA format



Sequence Data Preprocessing

This function wrapped **fastq-dump** function implemented in SRA Toolkit. See <http://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software> for details.

Input

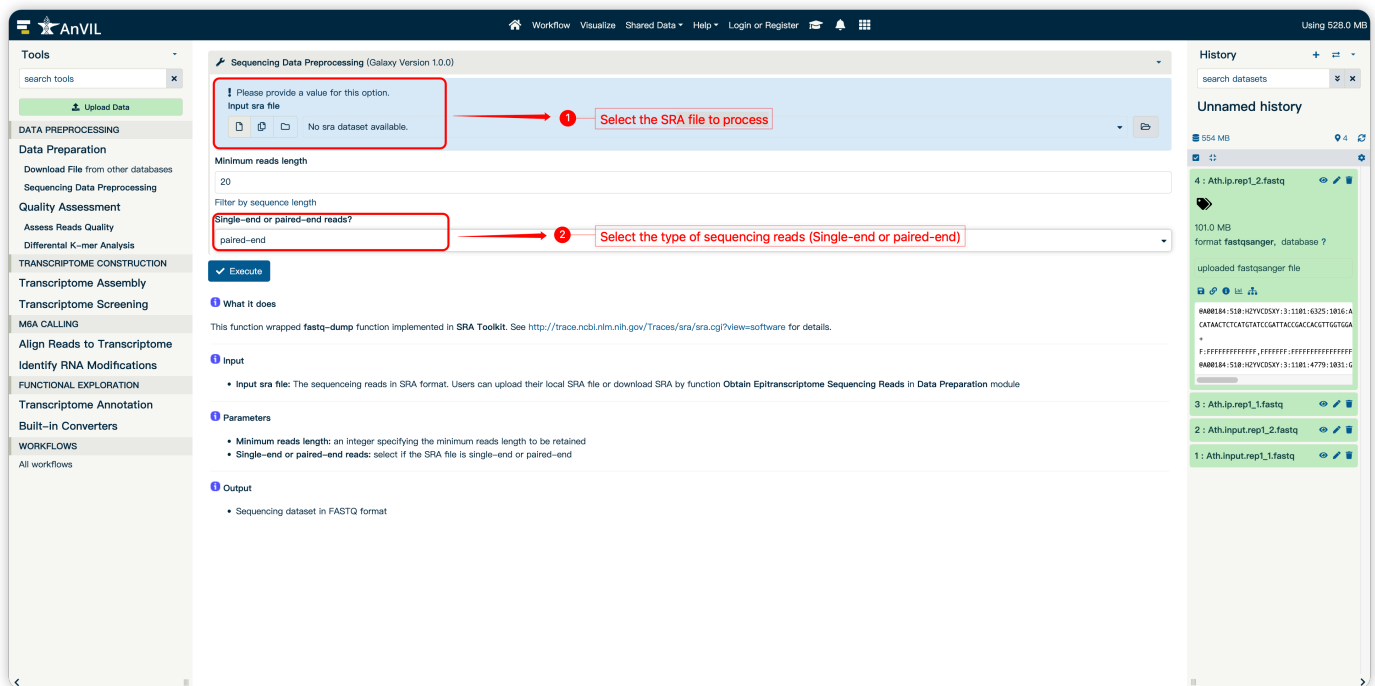
- **Input sra file:** The sequencing reads in SRA format. Users can upload their local SRA file or download SRA by function **Download File** in **Data Preparation** module

Output

- Sequencing dataset in FASTQ format

How to use this function

- The following screenshot shows us how to use this function to convert sequencing reads in SRA format to FASTQ format



Assess Reads Quality

In this function, two existing NGS tools MultiQC and fastp are integrated to check sequencing reads quality and obtain high-quality reads, respectively.

Input

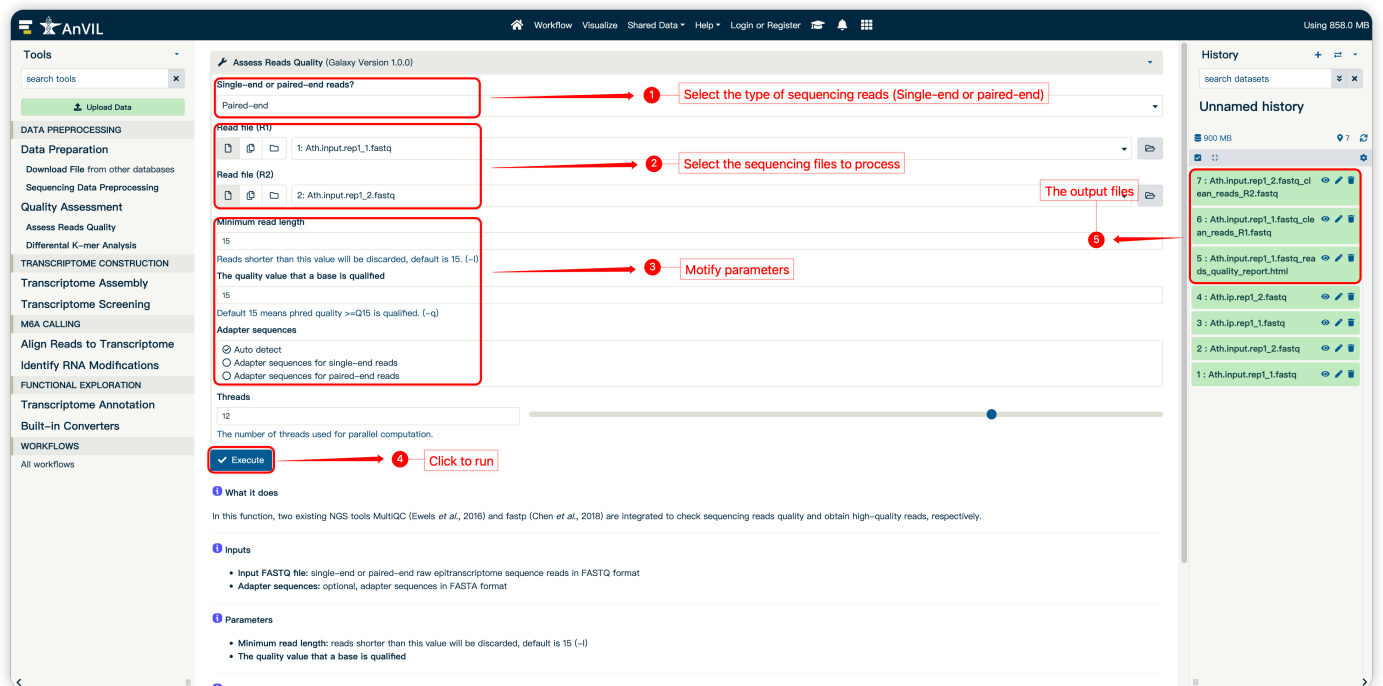
- **Input FASTQ file:** single-end or paired-end raw epitranscriptome sequence reads in FASTQ format
- **Adapter sequences:** optional, adapter sequences in FASTA format

Output

- Clean reads in FASTQ format
- Clean reads MultiQC report in HTML format

How to use this function

- The following screenshot shows us how to use this function to check m⁶A-Seq reads quality and obtain high-quality reads.



Differential *K*-mer Analyses

In this function, one existing disk-based program `kmdiff` is integrated to count *k*-mers from (possibly gzipped) FASTQ/FASTA files.

Input

- **Input FASTQ file:** IP and input cleaned epitranscriptome sequence reads in FASTQ format

Output

- **Differential *k*-mer sequences (input sample) in FASTA format**
- **Differential *k*-mer sequences (IP sample) in FASTA format**

How to use this function

- The following screenshot shows us how to use this function.

The screenshot displays the AnVIL web interface for the 'Differential K-mer Analysis' workflow. The interface is organized into three main sections: a left sidebar for navigation, a central workspace for workflow configuration, and a right sidebar for history and file management.

Left Sidebar (Navigation):

- Tools:** Includes a search bar and an 'Upload Data' button.
- DATA PREPROCESSING:**
 - Data Preparation:** Options for downloading files from databases, sequencing data preprocessing, and quality assessment.
 - Quality Assessment:** Options for assessing read quality and differential K-mer analysis.
- TRANSCRIPTOME CONSTRUCTION:**
 - Transcriptome Assembly:** The current workflow is selected here.
 - Transcriptome Screening:** Options for screening transcriptomes.
- M&A CALLING:** Options for aligning reads to transcriptomes, identifying RNA modifications, and functional exploration.
- Transcriptome Annotation:** Options for transcriptome annotation and built-in converters.
- WORKFLOWS:** A section for managing workflows.

Central Workspace (Workflow Configuration):

The workflow is titled 'Differential K-mer Analysis (Galaxy Version 1.0.0)'. It consists of several steps:

- Single-end or paired-end reads?**: A dropdown menu set to 'paired-end'.
- The FASTQ file(R1) or input**: A file selection box showing '6: Ath.input.rep1_1.fastq_clean_reads_R1.fastq'.
- The FASTQ file(R2) of input**: A file selection box showing '10: Ath.ip.rep1_2.fastq_clean_reads_R2.fastq'.
- The FASTQ file(R1) of IP**: A file selection box showing '9: Ath.ip.rep1_1.fastq_clean_reads_R1.fastq'.
- The FASTQ file(R2) of IP**: A file selection box showing '10: Ath.ip.rep1_2.fastq_clean_reads_R2.fastq'.
- k-mers length**: A text input field set to '8'.
- Threads**: A slider control set to '11'.

Below the workflow steps, there is a section for 'What it does' and 'Inputs'.

What it does:

In this function, one existing disk-based program kmrdiff is integrated to count k-mers from (possibly gzipped) FASTQ/FASTA files.

Inputs:

- Input FASTQ file: IP and Input cleaned epitranscriptome sequence reads in FASTQ format

Parameters:

- k-mers length: k-mer length (k from 8 to 20; default: 8)
- hard-min: min abundance to keep a k-mer

Right Sidebar (History and File Management):

The right sidebar contains a 'History' panel showing a list of previous runs. The most recent run (12) is highlighted, showing details such as 'Differential K-mer sequencing (IP sample)' and 'Differential K-mer sequencing (input sample)'. Below the history, there is a section for 'Unnamed history' and a list of files, including 'Ath.input.rep1_1.fastq_clean_reads_R1.fastq' and 'Ath.ip.rep1_2.fastq_clean_reads_R2.fastq'.