

# mlPEA User Manual

## (version 1.0)

- mlPEA is a user-friendly and multi-functionality platform specifically tailored to the needs of streamlined processing of m<sup>6</sup>A-Seq data in a reference genome-free manner. By taking advantage of machine learning (ML) algorithms, mlPEA enhanced the m<sup>6</sup>A-Seq data analysis by constructing robust computational models for identifying high-quality transcripts and high-confidence m<sup>6</sup>A-modified regions.
- mlPEA comprises four functional modules: **Data Preprocessing, Transcriptome Construction, m<sup>6</sup>A 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
<b>Download File</b>	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	<a href="#">SRA Toolkit</a>
<b>Sequence Data Preprocessing</b>	Convert epitranscriptome sequencing reads from SRA to FASTQ format	Epitranscriptome sequencing reads in SRA format	Epitranscriptome sequencing reads in FASTQ format	~2 mins	<a href="#">SRA Toolkit</a>
<b>Assess Reads Quality</b>	Check m <sup>6</sup> A-Seq reads quality and obtain high-quality reads	m <sup>6</sup> A-Seq reads in FASTQ format and adapter sequences in FASTA format	m <sup>6</sup> A-Seq reads in FASTQ format; MultiQC report in HTML	~2 mins	<a href="#">fastp</a> ; <a href="#">MultiQC</a>
<b>Differential K-mer Analyses</b>	Perform m <sup>6</sup> A-Seq differential k-mer analyses	m <sup>6</sup> A-Seq reads in FASTQ format	m <sup>6</sup> A-Seq reads in FASTA format	Depends on the file size	<a href="#">kmdiff</a>

## 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

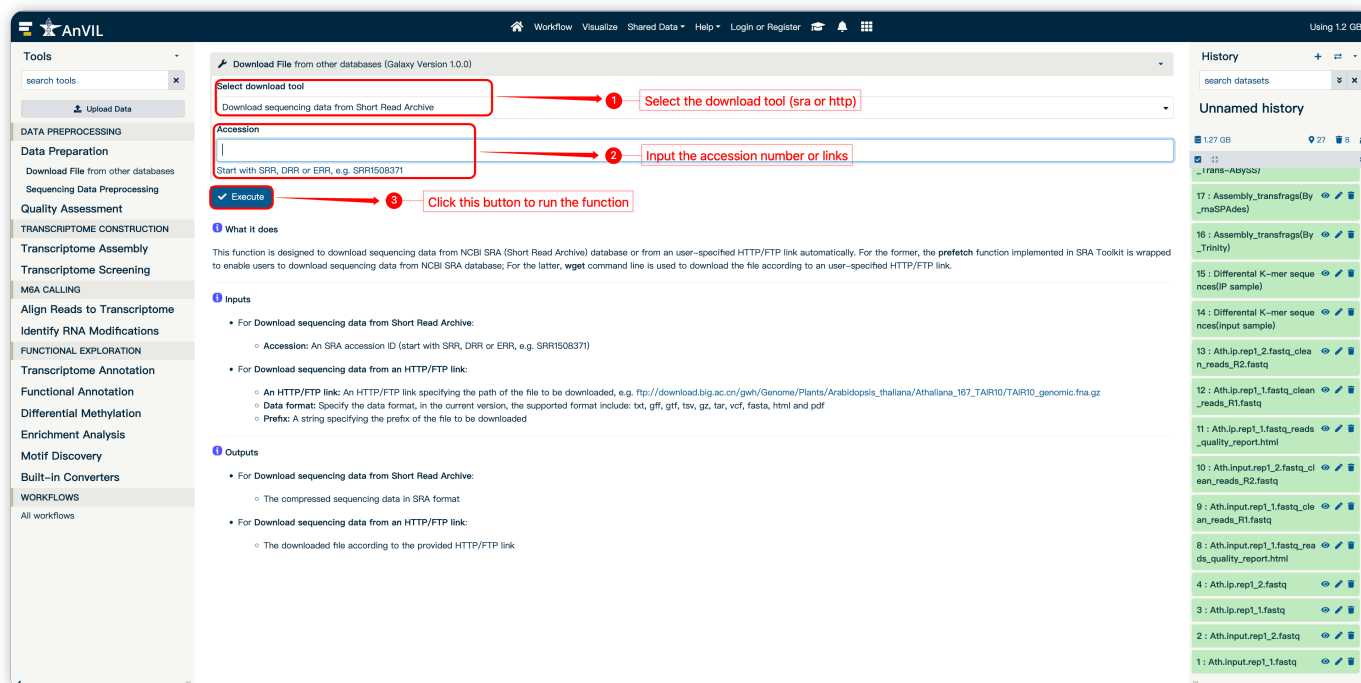
- For **Download sequencing data from Short Read Archive**:
  - **Accession**: An SRA accession ID (start with SRR, DRR or ERR, e.g. SRR1508371)
- For **Download sequencing data from an HTTP/FTP link**:
  - **An HTTP/FTP link**: An HTTP/FTP link specifying the path of the file to be downloaded, e.g. [ftp://download.big.ac.cn/gwh/Genome/Plants/Arabidopsis\\_thaliana/Athaliana\\_167\\_TAIR10/TAIR10\\_genomic.fna.gz](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

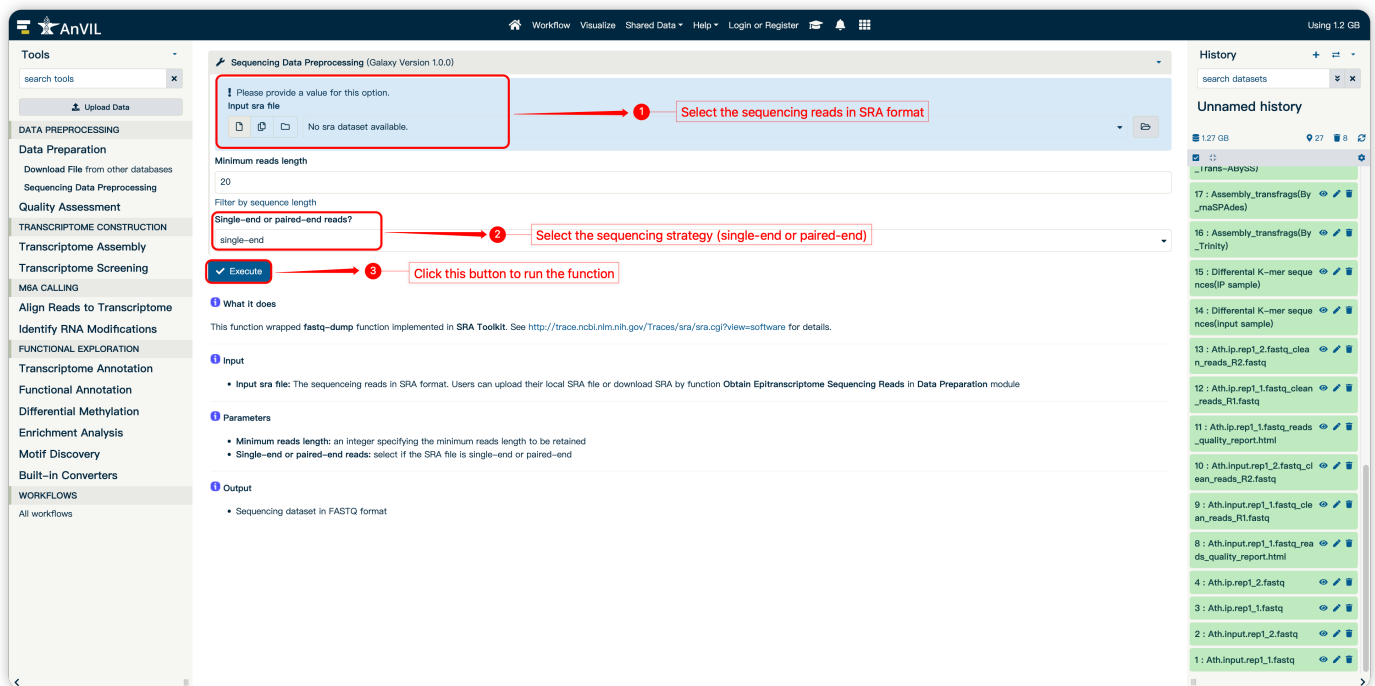
- **Input sra file:** The sequencing reads in SRA format. Users can upload their local SRA file or download SRA by function **Download File** in **Data Preparation** module

## 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

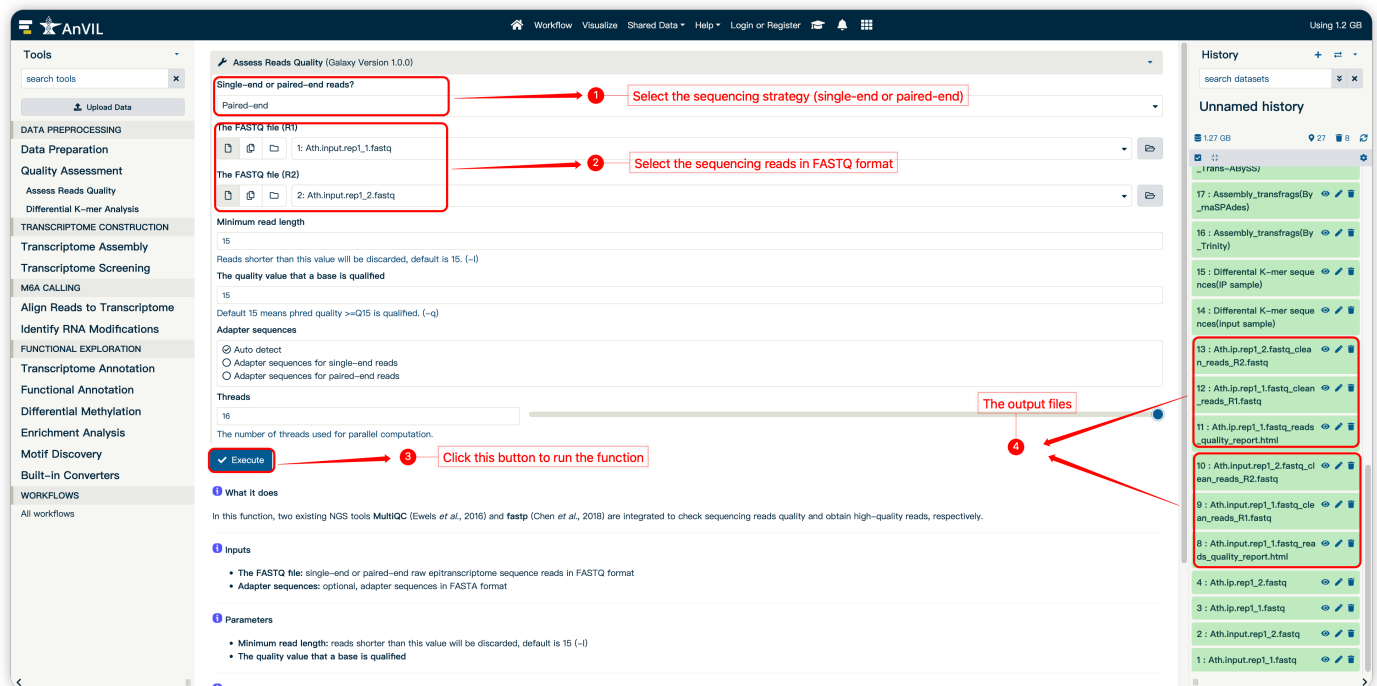
- **Input FASTQ file:** single-end or paired-end raw epitranscriptome sequence reads in FASTQ format
- **Adapter sequences:** optional, adapter sequences in FASTA format

## 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 m<sup>6</sup>A-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.

Tools

search tools

Upload Data

DATA PREPROCESSING

Data Preparation

Quality Assessment

Differential K-mer Analysis

TRANSCRIPTOME CONSTRUCTION

Transcriptome Assembly

Transcriptome Screening

M6A CALLING

Align Reads to Transcriptome

Identify RNA Modifications

FUNCTIONAL EXPLORATION

Transcriptome Annotation

Functional Annotation

Differential Methylation

Enrichment Analysis

Motif Discovery

Built-in Converters

WORKFLOWS

All workflows

Differential K-mer Analysis (Galaxy Version 1.0.0)

Single-end or paired-end reads?

paired-end

The FASTQ file(R1) of input

9: Ath.input.rep1\_1.fastq\_clean\_reads\_R1.fastq

The FASTQ file(R2) of input

13: Ath.ip.rep1\_2.fastq\_clean\_reads\_R2.fastq

The FASTQ file(R1) of IP

12: Ath.ip.rep1\_1.fastq\_clean\_reads\_R1.fastq

The FASTQ file(R2) of IP

13: Ath.ip.rep1\_2.fastq\_clean\_reads\_R2.fastq

k-mers length

8

k-mer length (k from 8 to 20; default: 8), (-k)

hard min

min abundance to keep a k-mer. (---hard-min)

20

Threads

12

The number of threads used for parallel computation.

Execute

What it does

In this function, one existing disk-based program **kmdiff** is integrated to count k-mers from (possibly gripped) 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

History

search datasets

Unnamed history

1: Ath.input.rep1\_1.fastq

2: Ath.input.rep1\_2.fastq

3: Ath.ip.rep1\_1.fastq

4: Ath.ip.rep1\_2.fastq

8: Ath.ip.rep1\_1.fastq\_reads\_report.html

9: Ath.input.rep1\_1.fastq\_clean\_reads\_R1.fastq

10: Ath.input.rep1\_2.fastq\_clean\_reads\_R2.fastq

11: Ath.ip.rep1\_1.fastq\_reads\_report.html

12: Ath.ip.rep1\_1.fastq\_clean\_reads\_R1.fastq

13: Ath.ip.rep1\_2.fastq\_clean\_reads\_R2.fastq

14: Differential K-mer sequence (input sample)

15: Differential K-mer sequence (IP sample)

16: Assembly\_transfrags (By Trinity)

17: Assembly\_transfrags (By masSPAdes)

18: Assembly\_transfrags (By Trinity)

1 Select the sequencing strategy(single-end or paired-end)

2 Select the cleaned reads of input sample in FASTQ format

3 Select the cleaned reads of IP sample in FASTQ format

4 Click this button to run the function

5 The output files