

# mlPEA User Manual

## (version 1.0)

- mlPEA is a user-friendly, full-functionality pipeline specifically designed to enhance the processing, analysis, and interpretation of m6A-Seq in non-model plants by leveraging machine learning.
- 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>
- mlPEA server can be accessed via <http://mlpea.omstudio.cloud>

## m<sup>6</sup>A Calling Module

This module provides step-by-step functions required for epitranscriptome reads mapping and identification of RNA modifications.

### Align Reads to Genome

Bowtie2 is wrapped to align epitranscriptome reads to genome.

Tools	Description	Input	Output	Time (test data)	Reference
<b>Bowtie2</b>	Bowtie2 is a short read aligner which achieves a combination of high speed, sensitivity and accuracy by combining the strengths of the full-text minute index with the flexibility and speed of hardware-accelerated dynamic programming algorithms, therefore bowtie2 is suitable for large genomes	Epitranscriptome sequencing reads in FASTQ format and reference genome sequences in FASTA format	Read alignments in SAM/BAM format	~5 mins	(Grabherr <i>et al.</i> , 2011)

### Identify RNA Modifications

**Identify RNA Modifications** implements three pipelines.

Tools	Description	Input	Output	Time (test data)	Reference
Peak Calling	used the SlidingWindow method with default options. To call m6A peaks, each transfrag was scanned using a 100-bp sliding window with 25-bp steps.	Input the transcripts in FASTA format;Input the ReadsCoverage file in IP sample;Input the ReadsCoverage file in input sample	The enriched peak region matrix in BED format	~5 mins	(Zhai <i>et al</i> , 2018)
Merge Biological Replicates	obtain consistent RNA modifications among three biological replicates.	Peak regions for biological Replicates	Consistent peak regions among multiple biological replicates	~10 mins	in-house scripts
High Confidence Peak	All peaks derived from the HC transcripts are designated as positive samples and utilized as input for the weakly supervised learning framework to facilitate model training.	Transcripts in FASTA format;The enriched peak region matrix in BED format	High confidence Peak region in BED format	~20mins	(Huang <i>et al</i> , 2021; Song <i>et al</i> , 2024)

## Align Reads to Genome

### Input

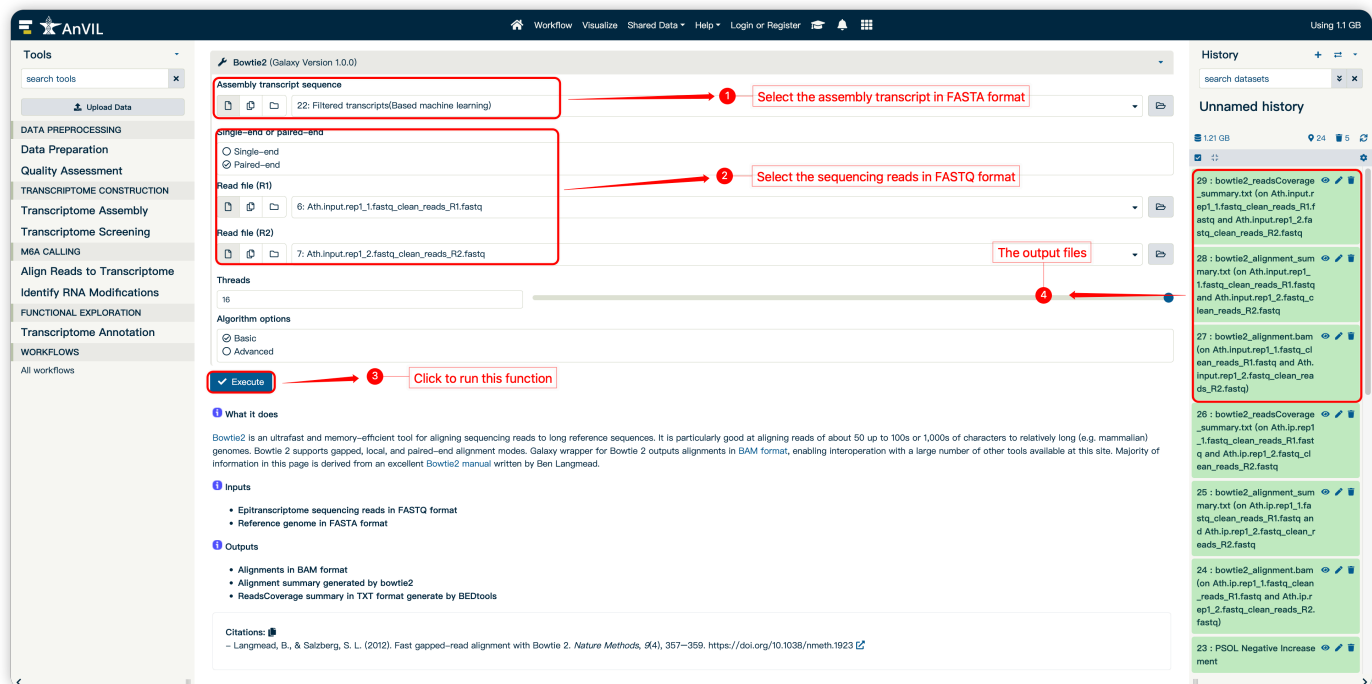
- Epitranscriptome sequencing reads in FASTQ format
- Assembled transcripts in FASTA format

### Output

- Alignments in BAM format
- Alignment summary generated by bowtie2
- ReadsCoverage summary in TXT format generate by BEDtools

### How to use this function

- The following screenshot shows us how to use this function



## Peak Calling

In this function, Peak calling was performed using PEA (Zhai *et al*, 2018), which used the SlidingWindow method with default options. To call m6A peaks, each transfrag was scanned using a 100-bp sliding window with 25-bp steps.

### Input

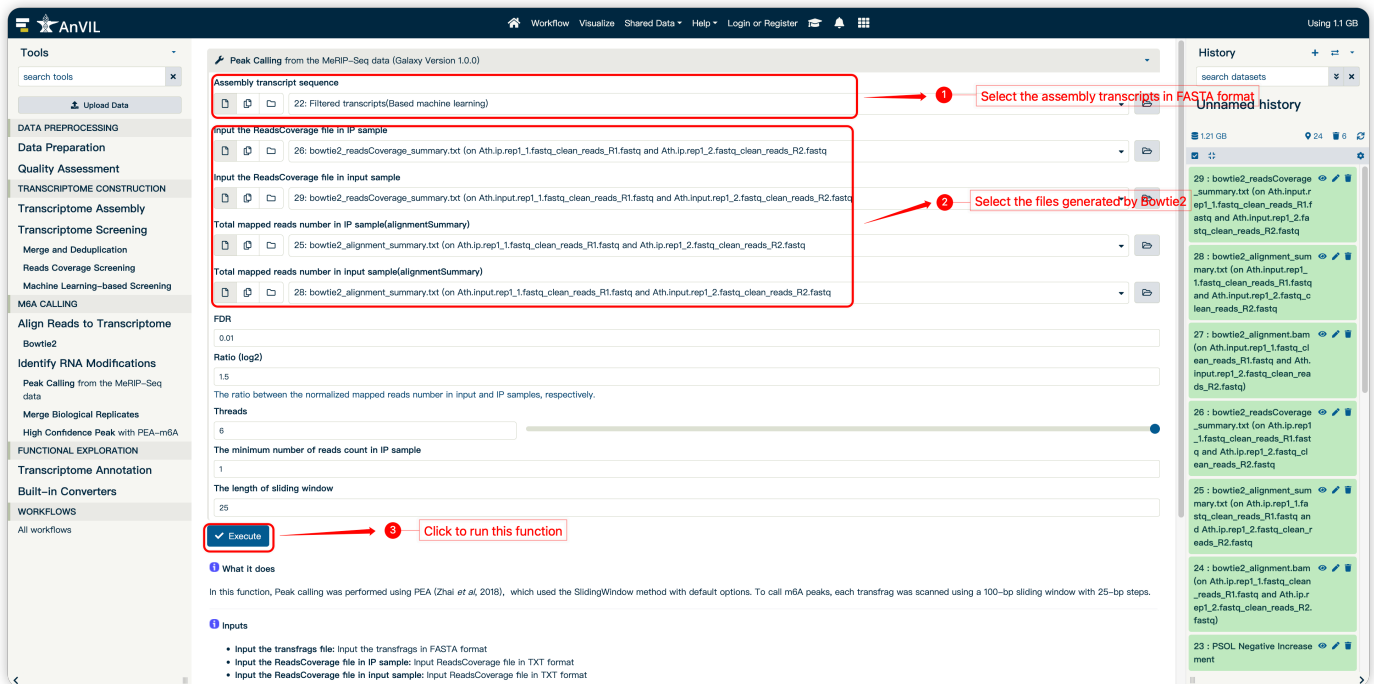
- **Input the transfrags file:** Input the transfrags in FASTA format
- **Input the ReadsCoverage file in IP sample:** Input ReadsCoverage file in TXT format
- **Input the ReadsCoverage file in input sample:** Input ReadsCoverage file in TXT format

### Output

- The enriched peak region matrix in BED format

### How to use this function

- The following screenshot shows us how to use this function



## Merge Biological Replicates

In this function, mlPEA obtained consistent RNA modifications among three biological replicates.

### Input

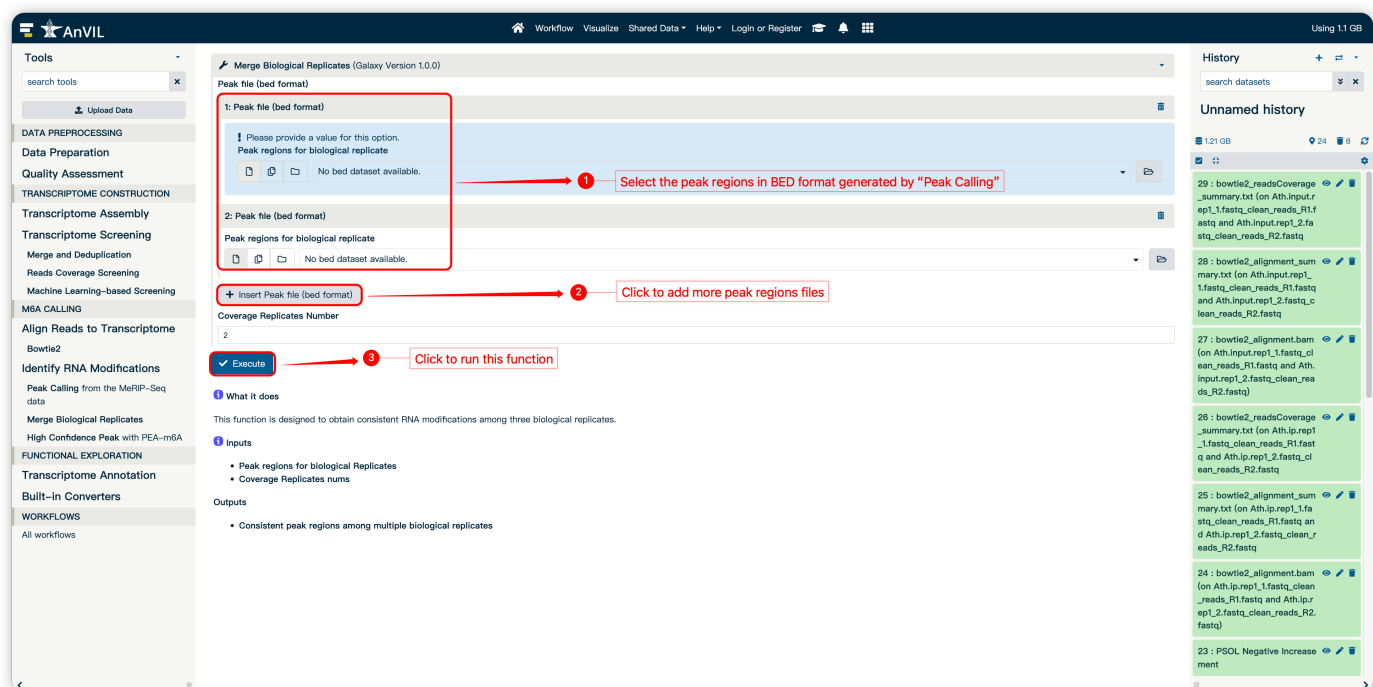
- Peak regions for biological Replicates

### Output

- Consistent peak regions among multiple biological replicates

### How to use this function

- The following screenshot shows us how to use this function



## High Confidence Peak

All peaks derived from the HC transcripts are designated as positive samples and utilized as input for the weakly supervised learning framework to facilitate model training.

### Input

- **Input the transfrags file:** Input the transfrags in FASTA format
- **Input the peak region file:** The enriched peak region matrix in BED format

### Output

- **High confidence Peak region in BED 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 'High Confidence Peak with PEA-m6A' workflow. The interface is structured into three main sections: a left sidebar for navigation, a central workspace for workflow execution, and a right sidebar for history and file management.

**Left Sidebar (Navigation):**

- Tools:** Includes a search bar and an 'Upload Data' button.
- DATA PREPROCESSING:** A section for data preparation.
- Data Preparation:** Includes 'Quality Assessment' and 'TRANSCRIPTOME CONSTRUCTION'.
- Transcriptome Assembly:** Includes 'Transcriptome Screening'.
- Transcriptome Screening:** Includes 'Merge and Deduplication', 'Reads Coverage Screening', and 'Machine Learning-based Screening'.
- M6A CALLING:** The current active section, containing:
  - Align Reads to Transcriptome
  - Bowtie2
  - Identify RNA Modifications
  - Peak Calling from the MeRIP-Seq data
  - Merge Biological Replicates
  - High Confidence Peak with PEA-m6A
  - FUNCTIONAL EXPLORATION**
  - Transcriptome Annotation
  - Built-in Converters
  - WORKFLOWS
  - All workflows

**Central Workspace:**

The main area displays the 'High Confidence Peak with PEA-m6A (Galaxy Version 1.0.0)' workflow. It includes a description, input fields, and a 'Run' button.

- Description:** 'Assembly transcript sequence' and 'Input the peak file in BED format'.
- Inputs:**
  - 22-Filtered transcripts(Based machine learning)
  - No bed dataset available.
- Parameters:**
  - Epoch: 50
  - Learning Rate: 0.0001
  - Initial learning rate: 1e-05
  - Decayed learning rate: The fixed-length sliding window length
  - The stride: 20
  - The steps length of sliding window: 50
  - The threshold of filtering HC peak: 0.5
- Run Button:** A green button with a checkmark and the text 'Execute'.

**Right Sidebar (History):**

The right sidebar shows a list of recent history items, including various transcript and alignment files.

- 29 : bowtie2\_readsCoverage\_summary.txt (on Ath.input.r
- 28 : bowtie2\_alignment\_summary.txt (on Ath.input.rep1
- 27 : bowtie2\_alignment.bam (on Ath.input.rep1.1.fastq.cl
- 26 : bowtie2\_readsCoverage\_summary.txt (on Ath.ip.rep1
- 25 : bowtie2\_alignment\_summary.txt (on Ath.ip.rep1.1.f
- 24 : bowtie2\_alignment.bam (on Ath.ip.rep1.1.fastq.clean
- 23 : PSOL Negative Increase ment

**Annotations:**

Three red boxes and arrows highlight key actions:

- 1:** Select the assembly transcript file in FASTA format
- 2:** Select the peak region file in BED format
- 3:** Click to run this function