

# 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>

## 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 Assembled Transcriptome

Bowtie2 is wrapped to align epitranscriptome reads to assembled transcriptome.

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 assembled transcriptome sequences in FASTA format	Read alignments in SAM/BAM format	~15 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	Using the SlidingWindow algorithm to identify regions with statistically significant enrichment of m <sup>6</sup> A signals compared to the background.	Assembled transcripts in FASTA format; Reads coverage file in IP and input sample; Total mapped reads number in IP and input sample	The enriched peak region matrix in BED format	~10 mins	(Zhai <i>et al</i> , 2018)
Merge Biological Replicates	Obtain consistent RNA modifications among multiple biological replicates.	Peak regions for biological Replicates	Consistent peak regions among multiple biological replicates in BED format	~1 mins	in-house scripts
ML-based peak Screening	All peaks derived from the assembled transcripts are designated as positive samples and utilized as input for the weakly supervised learning framework to facilitate model training.	Assembled transcripts in FASTA format; The enriched peak region matrix in BED format	High confidence Peak region in BED format	~6mins	(Huang <i>et al</i> , 2021; Song <i>et al</i> , 2024)

## Align Reads to Assembled Transcriptome

In this function, clean reads are aligned to high-quality transcripts using **Bowtie2**, with the *de novo* assembled transcriptome serving as the reference.

### Input

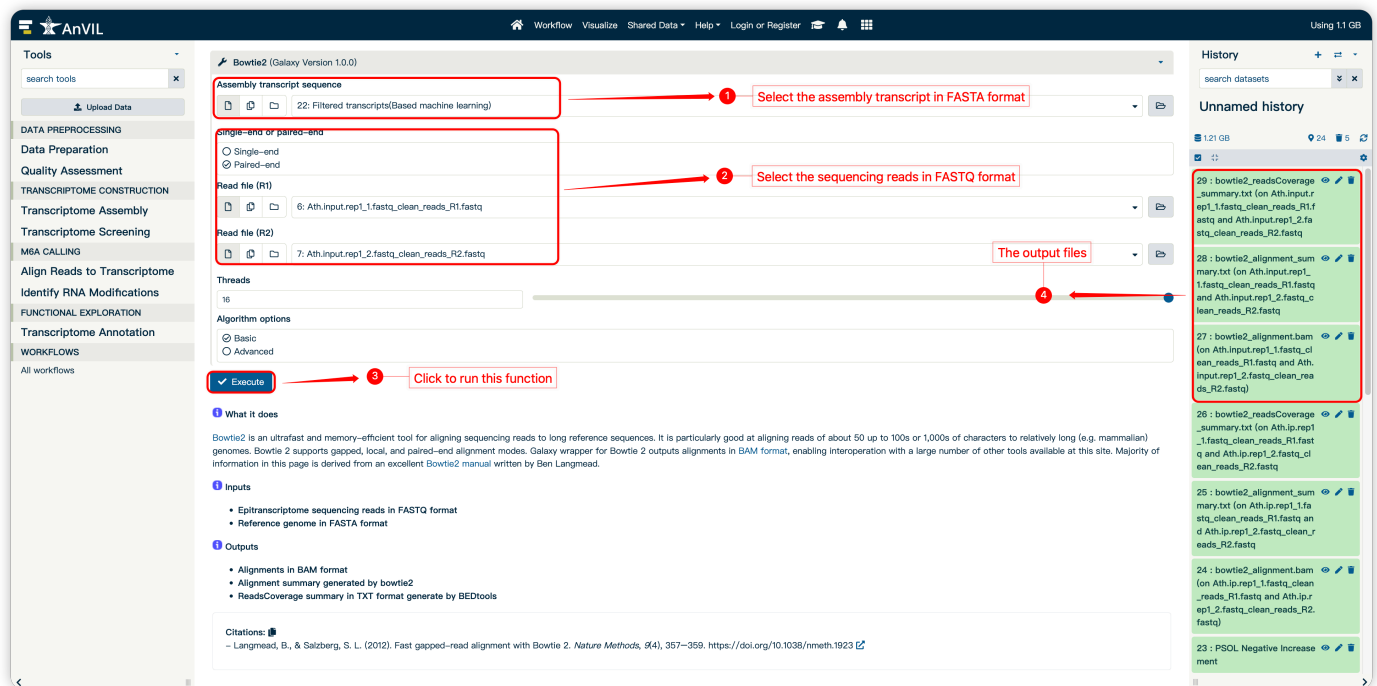
- Epitranscriptome sequencing reads in FASTQ format
- Assembled transcripts in FASTA format

### Output

- Alignment results in BAM format
- Alignment summary generated by bowtie2
- Reads coverage 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**, which used the SlidingWindow algorithm to identify regions with statistically significant enrichment of m<sup>6</sup>A signals compared to the background.

### Input

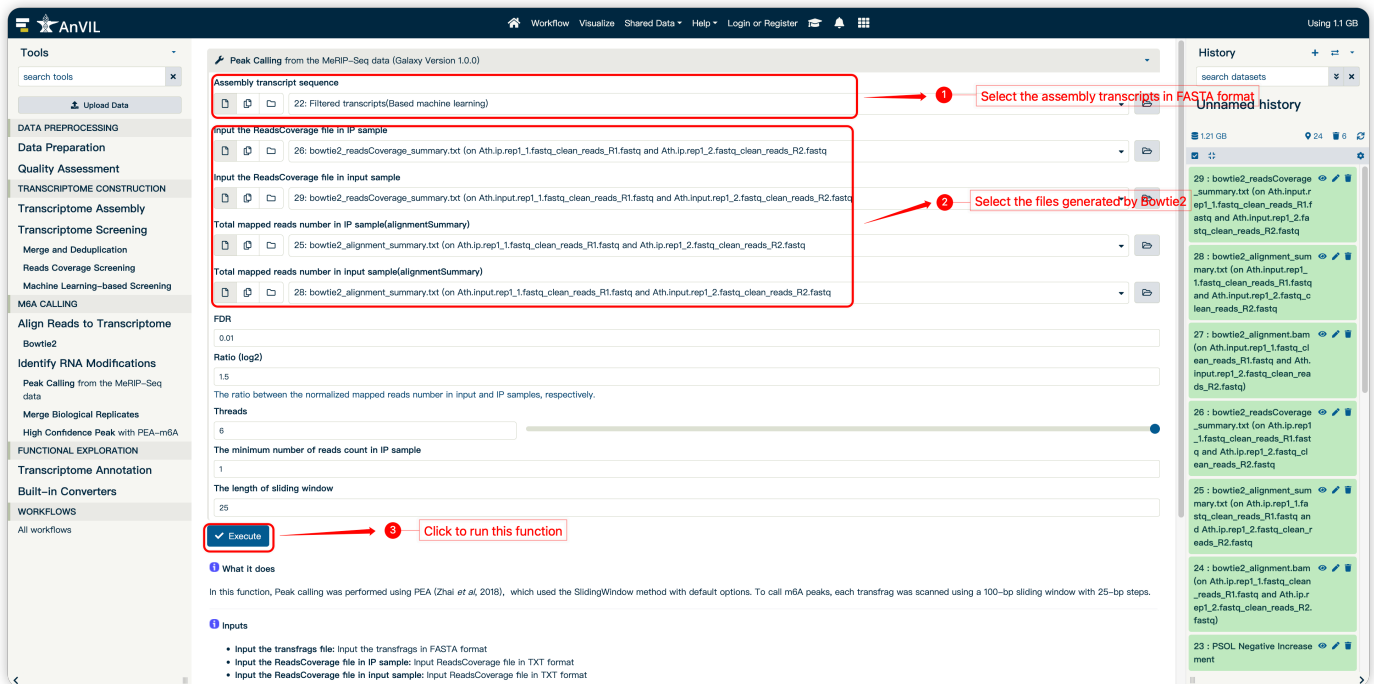
- **Input the assembled transcripts file:** assembled transcripts in FASTA format
- **Input the reads coverage files:** Reads coverage file in IP and input sample in TXT format
- **Input the alignment summary file:** Total mapped reads number in IP and input sample in TXT format

### 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 multiple biological replicates.

### Input

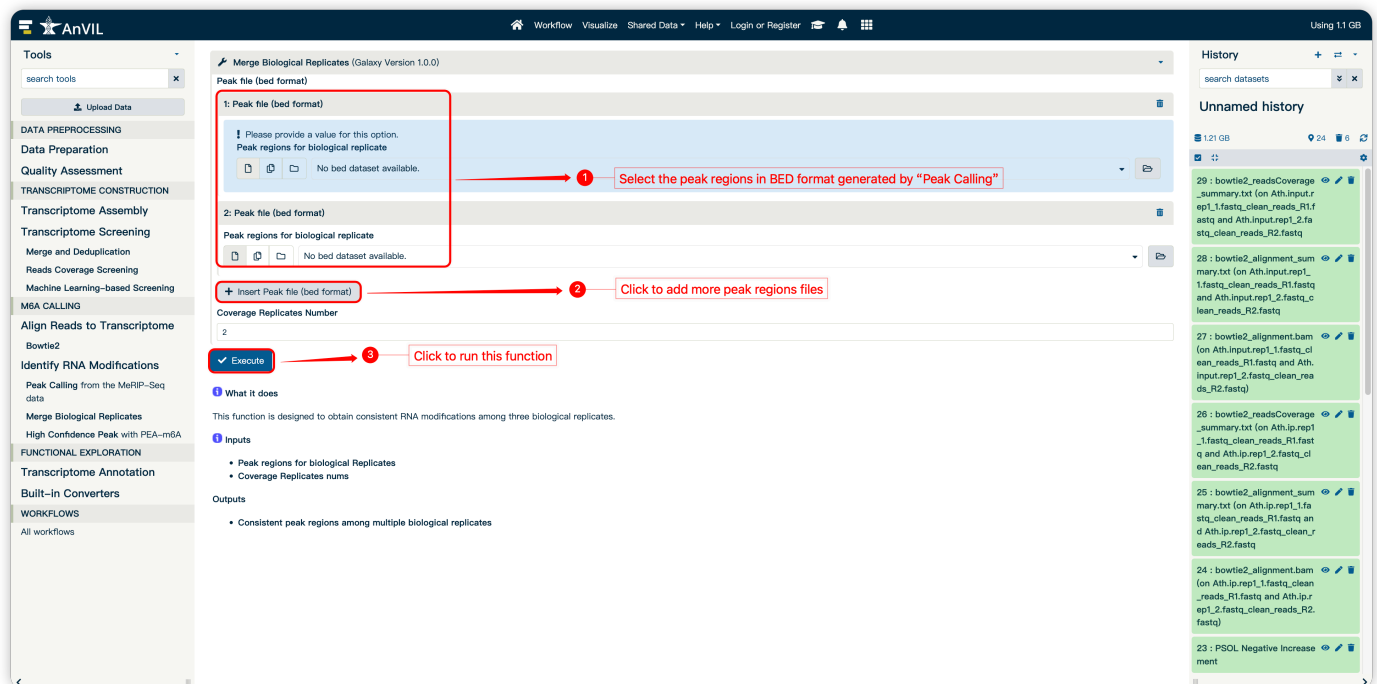
- Peak regions for biological replicates

### Output

- Consistent peak regions among multiple biological replicates in BED format

### How to use this function

- The following screenshot shows us how to use this function



## ML-based peak Screening

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

### Input

- **Input the assembled transcripts file:** assembled transcripts 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 organized into three main sections: a left sidebar for navigation, a central workspace for workflow execution, and a right sidebar for history and recent workflows.

**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:** Includes 'Align Reads to Transcriptome' and 'Bowtie2'.
- Identify RNA Modifications:** Includes 'Peak Calling from the MeRIP-Seq data'.
- FUNCTIONAL EXPLORATION:** Includes 'Merge Biological Replicates', 'High Confidence Peak with PEA-m6A', and 'Transcriptome Annotation'.
- Built-in Converters:** Includes 'WORKFLOWS' and 'All workflows'.

**Central Workspace:**

The central workspace displays the 'High Confidence Peak with PEA-m6A (Galaxy Version 1.0.0)' workflow. The workflow steps are as follows:

- Assembly transcript sequence:** A red box highlights the '22-Filtered transcripts(Based machine learning)' file. A red arrow points to this box with the label '1 Select the assembly transcript file in FASTA format'.
- Input the peak file in BED format:** A red box highlights the 'No bed dataset available.' message. A red arrow points to this box with the label '2 Select the peak region file in BED format'.
- Execute:** A red box highlights the 'Execute' button. A red arrow points to this box with the label '3 Click to run this function'.

**Right Sidebar (History):**

The right sidebar shows a 'History' panel with a search bar and a list of recent workflows. The workflows listed are:

- 29 : bowtie2\_readsCoverage
- 28 : bowtie2\_alignment\_summary.txt
- 27 : bowtie2\_alignment.bam
- 26 : bowtie2\_readsCoverage
- 25 : bowtie2\_alignment\_summary.txt
- 24 : bowtie2\_alignment.bam
- 23 : PSOL Negative Increase