

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>

Transcriptome Construction Module

This module provides step-by-step functions required for transcriptome construction.

Transcriptome Assembly

Four assemblers that support strand-specific m⁶A-Seq data are wrapped to use. Currently, **Trinity**, **rnaSPAdes**, **TransABYSS**, **TransLiG**.

Tools	Description	Input	Output	Time (test data)	Reference
Trinity	Trinity partitions the sequence data into many individual de Bruijn graphs, each representing the transcriptional complexity at a given gene or locus, and then processes each graph independently to extract full-length splicing isoforms and to tease apart transcripts derived from paralogous genes.	Sequencing reads in FASTQ format	Transcripts in FASTA format	~4 mins	(Grabherr <i>et al.</i> , 2011)
rnaSPAdes	rnaSPAdes has been developed on top of the SPAdes genome assembler and explores computational parallels between assembly of transcriptomes and single-cell genomes.	Sequencing reads in FASTQ format	Transcripts in FASTA format	~3 mins	(Prjibelski <i>et al.</i> , 2020)
TransABYSS	A <i>de novo</i> short-read transcriptome assembly and analysis pipeline that addresses variation in local read densities by assembling read substrings with varying stringencies and then merging the resulting contigs before analysis.	Sequencing reads in FASTQ format	Transcripts in FASTA format	~5 mins	(Robertson <i>et al.</i> , 2010)
TransLiG	TransLiG is shown to be significantly superior to all the salient <i>de novo</i> assemblers in both accuracy and computing resources when tested on artificial and real RNA-seq data.	Sequencing reads in FASTQ format	Transcripts in FASTA format	~4 mins	(Liu <i>et al.</i> , 2019)

Transcriptome Screening

Transcriptome screening implements three pipelines for transcripts deduplication, respectively.

Tools	Description	Input	Output	Time (test data)	Reference
Merge and Deduplication	A novel program CD-HIT for clustering biological sequences to reduce sequence redundancy and improve the performance of other sequence analyses.	Multiple Transcripts files in FASTA format	Transcripts in FASTA format	~5 mins	(Fu <i>et al.</i> , 2012)
Reads Coverage-based Screening	An ultrafast and memory-efficient tool Bowtie 2 for aligning sequencing reads, reads coverage was measured by SAMtools and BEDTools.	Transcripts in FASTA format	Transcripts in FASTA format	~20 mins	(Langmead <i>et al.</i> , 2009; Langmead <i>et al.</i> , 2012) (Danecek <i>et al.</i> , 2021) (Quinlan <i>et al.</i> , 2010)
Machine Learning-based Screening	An ML-based classification model to distinguish high-quality transcripts from noisy ones using the random forest-based PSoL algorithm, which requires only a set of positive samples. The positive sample set is constructed by clustering analysis of assembled transcripts and already annotated mRNAs using a fast and sensitive sequence alignment method MMseqs2. It provides more than 670 features to encode each transcript sequence.	Transcripts in FASTA format	Transcripts in FASTA format	~10mins	(Steinegger <i>et al.</i> , 2017) (Wang <i>et al.</i> , 2023) (Wang <i>et al.</i> , 2006)

Transcriptome Assembly

Currently, mlPEA wrapped four assemblers that support strand-specific m⁶A-Seq data to *de novo* assemble transcripts. Here, we take Trinity as an example to show how to use mlPEA to run transcriptome assembly, the other three assemblers are similar.

Input

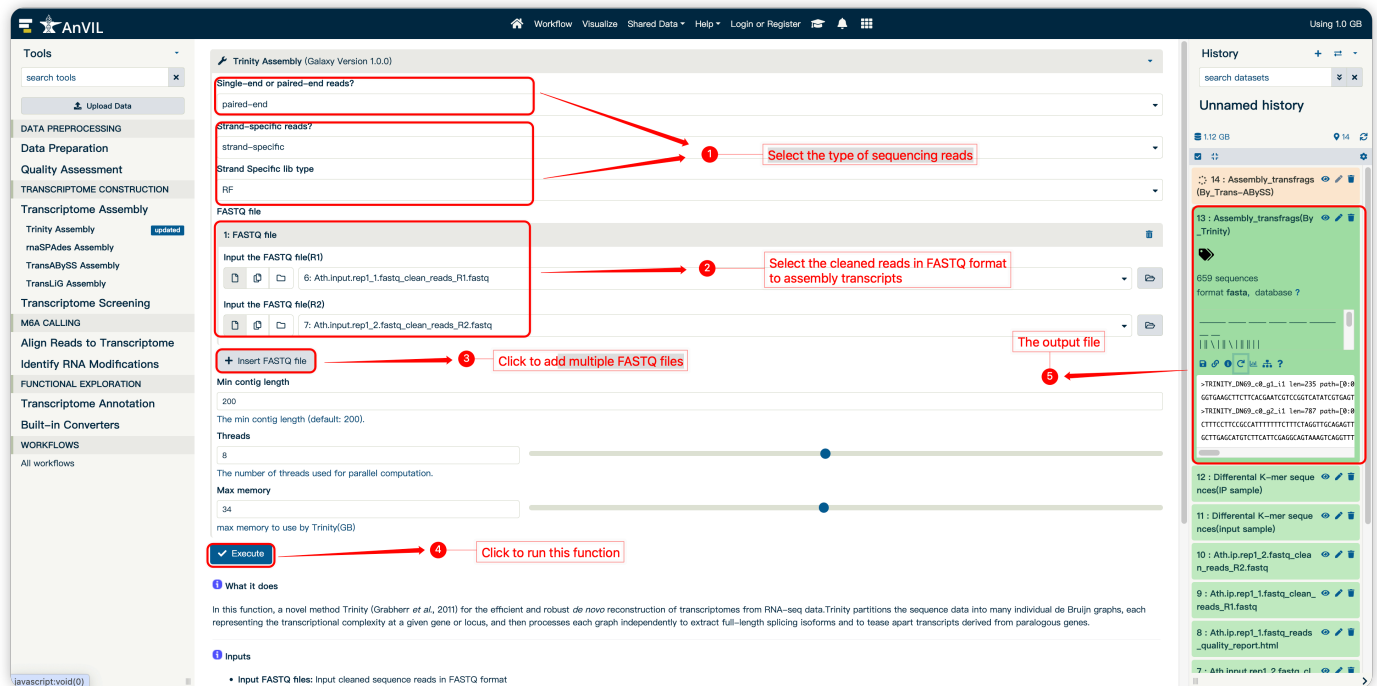
- **Input FASTQ files:** cleaned m⁶A-Seq input reads in FASTQ format

Output

- *De novo* assembled transcripts by Trinity in FASTA format

How to use this function

- The following screenshot shows us how to use this function



Merge and Deduplication

In this function, a novel program **CD-HIT** for clustering biological sequences to reduce sequence redundancy and improve the performance of other sequence analyses.

Input

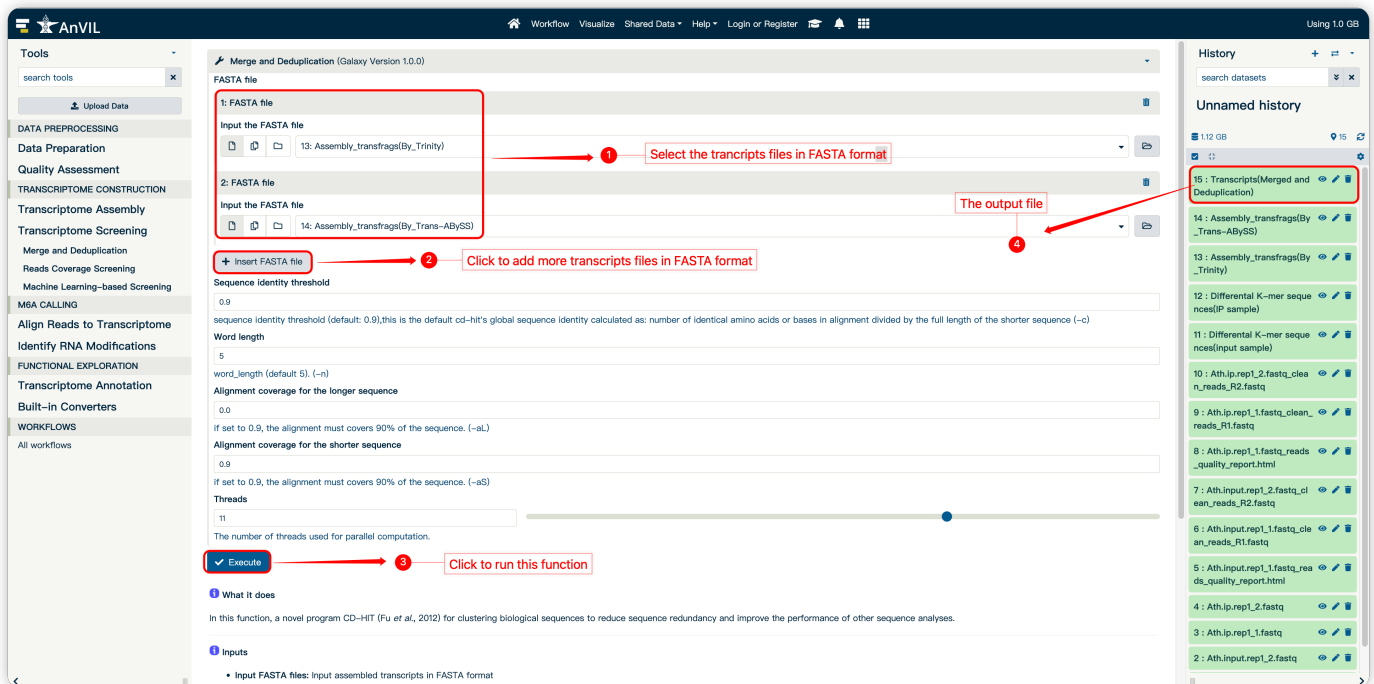
- **Input FASTA files:** assembled transcripts in FASTA format

Output

- **Transcripts after deduplication in FASTA format**

How to use this function

- The following screenshot shows us how to use this function



Reads Coverage-based Screening

In this function, we developed a pipeline for screening transcripts based on reads coverage. **Bowtie 2** for aligning sequencing reads, reads coverage was measured by **SAMtools** and **BEDTools**.

Input

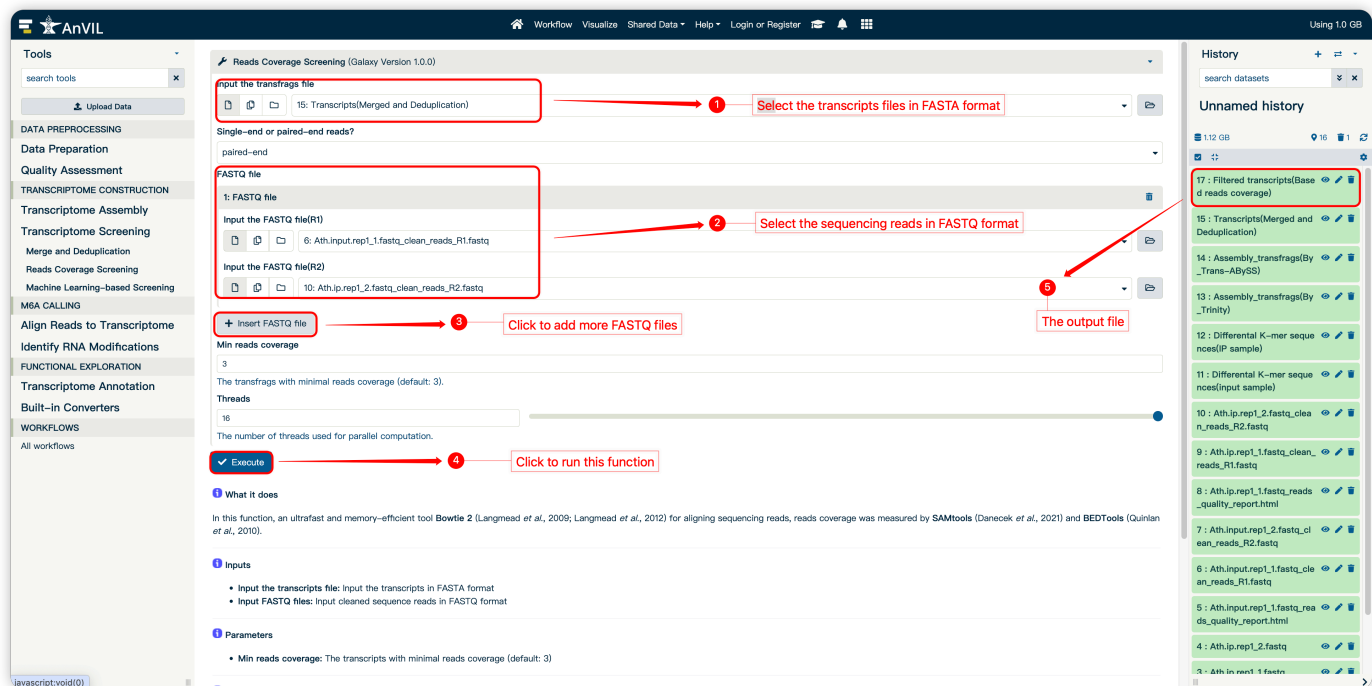
- **Input FASTQ files:** cleaned m⁶A-Seq input reads in FASTQ format

Output

- **De novo** assembled transcripts in FASTA format

How to use this function

- The following screenshot shows us how to use this function



Machine Learning-based Screening

In this function, an ML-based classification model to distinguish high-quality transcripts from noisy ones using the random forest-based Positive Sample only Learning (PSoL) algorithm, which requires only a set of positive samples. The positive sample set is constructed by clustering analysis of assembled transcripts and already annotated mRNAs using a fast and sensitive sequence alignment method **MMseqs2**. It provides more than 670 features to encode each transcript sequence, including 177 sequence-intrinsic features, 399 physico-chemical features and 101 structure-based features.

Input

- **Input FASTQ files:** *de novo* assembled transcripts in FASTA format

Output

- **High-quality transcripts in FASTA format**
- **PSoL Negative Increase in PDF format**
- **RNA features based on corain package in CSV 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

TRANSCRIPTOME CONSTRUCTION

Transcriptome Assembly

Transcriptome Screening

Merge and Deduplication

Reads Coverage Screening

Machine Learning-based Screening

MSA CALLING

Align Reads to Transcriptome

Identify RNA Modifications

FUNCTIONAL EXPLORATION

Transcriptome Annotation

Functional Annotation

Enrichment Analysis

Motif Discovery

Built-in Converters

WORKFLOWS

All workflows

Machine Learning-based Screening (Galaxy Version 1.0.0)

Input the transfrags file

15: Transcripts(Merged and Deduplication)

1 Select the transcripts in FASTA format

Positive Sample Generation

2 Modify the parameters

Feature Encoding

Model Learning and Labeling

Whether to perform k-fold cross-validation?

Yes

No

k-fold cross validation

5

Threads

10

The number of threads used for parallel computation.

3 Click to run this function

4 The output file

What it does

In this function, we utilized MMSeqs2 (Steinegger *et al.*, 2017) to cluster the assembled transfrags with both mRNA and ncRNA sequences from the Ensembl Plants across all plant species. Both assembled transfrags and annotated plant sequences to serve as positive sample set. Then we used the corain package (Wang *et al.*, 2023) to encode for RNA from three perspectives, sequence, structure and physical chemical properties. Adopted the improved PSOL(IPSO) algorithm (Wang *et al.*, 2006), the initial negative samples were generated from "unlabeled" sample set based on the PU bagging algorithm, the expanded negative samples iteratively using RF classifier until the designated iteration number was reached.

Inputs

Input the transcripts file: Input the transcripts in FASTA format

Parameters

Min reads coverage: The transcripts with minimal reads coverage (default: 3)

Outputs

Filtered transcripts with machine learning in FASTA format

PSOL Negative Increaseement in PDF format

RNA features based on corain package in CSV format

History

search datasets

Unnamed history

1.12 GB

23: PSOL Negative Increaseement

22: Filtered transcripts (Based machine learning)

17: Filtered transcripts(Base d reads coverage)

15: Transcripts(Merged and Deduplication)

14: Assembly_transfrags(By _Trans-ABYSS)

13: Assembly_transfrags(By _Trinity)

12: Differential K-mer sequen ces(IP sample)

11: Differential K-mer sequen ces(input sample)

10: Ath.ip.rep1_2.fastq_clean_reads_R2.fastq

9: Ath.ip.rep1_1.fastq_clean_reads_R1.fastq

8: Ath.ip.rep1_1.fastq_reads _quality_report.html

7: Ath.input.rep1_2.fastq_clean_reads_R2.fastq

6: Ath.input.rep1_1.fastq_clean_reads_R1.fastq