

## **SmallRNA Method (Illumina)**

# **1. Experimental Procedure**

## **1.1 Sample Quality Control**

Please refer to QC report for methods of sample quality control.

## **1.2 Library Construction, Quality Control and Sequencing**

Total RNA was used as input material for the RNA sample preparations. Briefly, 3' and 5' adaptors were ligated to 3' and 5' end of small RNA, respectively. Then the first strand cDNA was synthesized after hybridization with reverse transcription primer. The double-stranded cDNA library was generated through PCR enrichment. After purification and size selection, libraries with insertions between 18~40 bp were ready for sequencing.

After the library construction is completed, quantification of the inserted fragments and the effective concentration to ensure the quality of the library.

The qualified libraries were pooled and sequenced on Illumina platforms, according to effective library concentration and data amount required.

# **2. Bioinformatics Analysis Pipeline**

## **2.1 Data Quality Control**

### **2.1.1 Raw Data**

The original fluorescence image files obtained from Illumina platform are transformed to short reads (Raw data) by base calling and these short reads are recorded in FASTQ format (Cock P. et al, 2010), which contains sequence information and corresponding sequencing quality information.

Raw data (raw reads) of fastq format were firstly processed through fastp software.

In this step, clean data (clean reads) were obtained by removing reads containing adapter, reads containing poly-N and low-quality reads from raw data. At the same time, Q20, Q30 and GC content of the clean data were calculated. All the downstream analyses were based on the clean data with high quality.

### **2.1.2 Evaluation of Data (Data Quality Control)**

Sequence artifacts, including reads containing adapter contamination, low-quality nucleotides and unrecognizable nucleotide (N), undoubtedly set the barrier for the subsequent reliable bioinformatics analysis. Hence quality control is an essential step and applied to guarantee the meaningful downstream analysis. we used Fastp (version 0.23.1) (Chen S. et al, 2018) to perform basic statistics on the quality of the raw reads.

The steps of data processing were as follows:

- (1) Discard a paired reads if either one read contains adapter contamination;
- (2) Discard a paired reads if more than 10% of bases are uncertain in either one read;
- (3) Discard a paired reads if the proportion of low quality (Phred quality <5) bases is over 50% in either one read.

### **2.1.3 Mapping to Reference Genome**

Bowtie (Langmead et al, 2009) is used to locate sRNA after length screening to the reference sequence, and analyzed the distribution of sRNA on the reference sequence.

### **2.1.4 Analysis of Known miRNA**

Reads matched to reference sequences are compared with sequences in the

specified range in miRBase. The sRNA information of each sample is obtained by mirdeep2 (Friedlander et al, 2011) and SRna-Tools-CLI. Including the known miRNA secondary structure, miRNA sequence, length, frequency of occurrence, etc.

### **2.1.5 Remove Reads from These Sources**

Annotate sRNA using the ncRNA sequences of the species, or select rRNA, tRNA, snRNA, and snoRNA from RFAM to annotate sRNA, identifying and removing potential rRNA, tRNA, snRNA, and snoRNA. Using species-specific repeat sequence annotation information, or reference sequence information for de novo prediction of repeat sequences, align sRNA with repeat sequences to identify and remove potential repeat sequences, and statistically analyze the types and quantities of sRNA matching various repeat types.

### **2.1.6 Analysis of NAT-siRNA (Plant)**

Natural antisense transcripts (NATs), as one type of regulatory RNAs, occur prevalently in plant genomes and play significant roles in physiological and pathological processes. NAT can be divided into two categories: cis-NAT and trans-NAT. NAT-siRNA is detected by NAT gene of corresponding species in PlantNATsDB database.

### **2.1.7 Novel miRNA Prediction**

The characteristic hairpin structure of miRNA precursors can be used to predict new miRNAs. Using miREvo (Wen et al., 2012) and mirdeep2 (Friedlander et al., 2011), sRNA sequences of a certain length are extracted and aligned with the reference genome. These sequences are then analyzed for secondary structure, Dicer cleavage

site information, energy characteristics, and other features to predict novel miRNAs. Statistically analyze the aligned sRNA sequences for their sequence, length, occurrence frequency, the distribution of the first nucleotide at different miRNA lengths, and the nucleotide distribution at each position for all miRNAs.

### **2.1.8 Analysis of Ta-siRNA (Plant)**

Trans-acting siRNAs (ta-siRNAs) are a new class of small RNAs that regulate plant development. The identification of known TAS gene based on Arabidopsis and Oryza sativa database. We used UEA sRNA tools (Moxon et al., 2008) to predict new TAS gene.

### **2.1.9 Small RNA Annotation**

Summarize the alignment and annotation of all small RNAs with various types of RNAs. Since a single sRNA can match multiple different annotation categories, to ensure that each unique small RNA is assigned to only one annotation, follow the priority order: known miRNA > rRNA > tRNA > snRNA > snoRNA > repeats > genes > NAT-siRNA > novel miRNA.

### **2.1.10 miRNA Base Edit**

MicroRNA may undergo nucleotide editing at certain positions, leading to changes in the seed sequence and subsequently altering target genes (Wei et al., 2009). By aligning the sRNA sequences from each sample with the detected known and novel mature miRNAs as well as their precursors, potentially mutated miRNAs can be identified.

### **2.1.11 miRNA Family Analysis**

Conduct a family analysis of the detected known and novel miRNAs to explore the presence of their miRNA families in other species. Known miRNAs can be identified using miFam.dat to determine their family origin, while novel miRNAs can be classified using RFAM to determine their RFAM family.

#### **2.1.12 miRNA Expression and Differential Expression**

Statistically analyze the expression levels of known and novel miRNAs in each sample, and normalize the expression levels using TPM (Zhou et al., 2010).  $TPM = (\text{read count} * 1,000,000) / \text{libsize}$  (libsize: total miRNA read count). For samples with biological replicates, use DESeq2 (Love MI et al., 2014) for differential expression analysis between two comparison groups. DESeq2 provides statistical routines to determine differential expression in digital gene expression data using a model based on the negative binomial distribution. For samples without biological replicates, use the edgeR (Robinson MD et al., 2010) TMM algorithm to normalize read count data for analysis.

#### **2.1.13 Target Gene Prediction for Known and Novel miRNA**

For animals, use miRanda and RNAhybrid to predict miRNA target genes, taking the intersection as the final targeting result. For plants, use psRobot to predict miRNA target genes. Perform target gene prediction for the analyzed known and novel miRNAs to obtain the relationships between miRNAs and their target genes.

#### **2.1.14 Enrichment Analysis**

We used clusterProfiler (Yu G et al, 2012) to achieve functional enrichment analysis of differentially expressed genes in GO (Gene Ontology). Differentially expressed

genes are significantly enriched, and  $\text{padj} < 0.05$  as the threshold of significant enrichment. KEGG (Kyoto Encyclopedia of Genes and Genomes) is the main public database of pathway significant enrichment analysis, hypergeometric test is applied to identify the pathway of significant enrichment in candidate target genes, and  $\text{padj} < 0.05$  as the threshold of significant enrichment. The differentially expressed genes in the KEGG pathway is analyzed by clusterProfiler.

### **3. References**

Cock P.J.A. et al (2010). The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants. *Nucleic acids research* 38, 1767-1771.

Shifu Chen, Yanqing Zhou, Yaru Chen, Jia Gu; fastp: an ultra-fast all-in-one FASTQ preprocessor, *Bioinformatics*, Volume 34, Issue 17, 1 September 2018, Pages i884–i890.

Langmead B, Trapnell C, Pop M, Salzberg SL. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol.* 2009;10(3):R25. doi:10.1186/gb-2009-10-3-r25. (Bowtie)

Friedländer MR, Mackowiak SD, Li N, Chen W, Rajewsky N. miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res.* 2012;40(1):37-52. doi:10.1093/nar/gkr688. (miRDeep2)

Wen M, Shen Y, Shi S, Tang T. miREvo: an integrative microRNA evolutionary analysis platform for next-generation sequencing experiments. *BMC Bioinformatics.* 2012;13:140. Published 2012 Jun 21. doi:10.1186/1471-2105-13-140.

Moxon S, Schwach F, Dalmay T, Maclean D, Studholme DJ, Moulton V. A toolkit for analysing large-scale plant small RNA datasets. *Bioinformatics*. 2008;24(19):2252-2253. doi:10.1093/bioinformatics/btn428.

Wei Y, Chen S, Yang P, Ma Z, Kang L. Characterization and comparative profiling of the small RNA transcriptomes in two phases of locust. *Genome Biol*. 2009;10(1):R6. doi:10.1186/gb-2009-10-1-r6.

Zhou L, Chen J, Li Z, et al. Integrated profiling of microRNAs and mRNAs: microRNAs located on Xq27.3 associate with clear cell renal cell carcinoma. *PLoS One*. 2010;5(12):e15224. Published 2010 Dec 30. doi:10.1371/journal.pone.0015224.

Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014;15(12):550. doi:10.1186/s13059-014-0550-8. (DESeq2)

Robinson MD, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*. 2010;26(1):139-140. doi:10.1093/bioinformatics/btp616.

Yu G, Wang LG, Han Y, He QY. clusterProfiler: an R package for comparing biological themes among gene clusters. *OMICS*. 2012;16(5):284-287. doi:10.1089/omi.2011.0118.