

## 1 **Methods**

### 2 3 1. RNA extension and transcriptome analysis RNA

#### 4 1.1 Sample collection and preparation

5 RNA degradation and contamination was monitored on 1% agarose gels. RNA  
6 purity was checked using the Nano Photometer® spectrophotometer (IMPLEN, CA,  
7 USA). RNA concentration was measured using Qubit® RNA Assay Kit in Qubit®2.0  
8 Fluorometer (Life Technologies, CA, USA). RNA integrity was assessed using the  
9 RNA Nano 6000 Assay Kit of the Agilent Bioanalyzer 2100 system (Agilent  
10 Technologies, CA, USA).

#### 11 12 1.2 cDNA Library preparation for sequencing

13 A total amount of 2 µg RNA per sample was used as input material for the RNA  
14 sample preparations. Sequencing libraries were generated using VAHTSTM  
15 mRNA-seq V2 Library Prep Kit for Illumina® following manufacturer's  
16 recommendations and index codes were added to attribute sequences to each sample.  
17 Briefly, mRNA was purified from total RNA using poly-T oligo-attached magnetic  
18 beads. Fragmentation was carried out using divalent cations under elevated  
19 temperature in VAHTSTM First Strand Synthesis Reaction Buffer (5X) . First strand  
20 cDNA was synthesized using random hexamer primer and M-MuLV Reverse  
21 Transcriptase ( RNase H- ) . Second strand cDNA synthesis was subsequently  
22 performed using DNA polymerase I and RNase H. Remaining overhangs were  
23 converted into blunt ends via exonuclease/polymerase activities. After adenylation of  
24 3' ends of DNA fragments, Adaptor were ligated to prepare for library. In order to  
25 select cDNA fragments of preferentially 150~200 bp in length, the library fragments  
26 were purified with AMPure XP system (Beckman Coulter, Beverly, USA). Then 3 µl  
27 USER Enzyme (NEB, USA) was used with size-selected, adaptor-ligated cDNA at  
28 37°C for 15 min followed by 5 min at 95 °C before PCR. Then PCR was performed  
29 with Phusion High-Fidelity DNA polymerase, Universal PCR primers and Index (X)  
30 Primer. At last, PCR products were purified (AMPure XP system) and library quality  
31 was assessed on the Agilent Bioanalyzer 2100 system. The libraries were then  
32 quantified and pooled. Paired-end sequencing of the library was performed on the  
33 HiSeq XTen sequencers (Illumina, San Diego, CA).

#### 34 35 1.3 Data assessment and quality control

36 **FastQC** (version 0.11.2) was used for evaluating the quality of sequenced data.  
37 Raw data was filtered by **Trimmomatic** (version 0.36) according to several steps: 1)  
38 Removing sequences with N bases; 2) Removing adaptor sequence if reads contains;  
39 3) Removing low quality bases from reads 3' to 5' (Q < 20); 4) Removing low quality  
40 bases from reads 5' to 3' (Q < 20); 5) Using a sliding window method to remove the  
41 base value less than 20 of reads tail (window size is 5 bp); 6) Removing reads with  
42 reads length less than 35nt and its pairing reads. And the remaining clean data was  
43 used for further analysis.

44

#### 45 1.4 Transcriptome assembly and gene annotation

46 The remaining clean reads were de novo assembled into transcripts using **Trinity**  
47 (version 2.0.6) with default settings. Transcripts with a minimum length of 200 bp  
48 were clustered to minimize redundancy. For each cluster (representing the  
49 transcriptional complexity for the same gene), the longest sequence was preserved and  
50 designated as unigene. Unigenes were blasted against NCBI Nr (NCBI non-redundant  
51 protein database), SwissProt, TrEMBL, CDD (Conserved Domain Database), Pfam  
52 and KOG (eukaryotic Orthologous Groups) databases (E-value < 1e-5). According to  
53 the priority order of the best aligned results of NR, SwissProt and TrEMBL to  
54 determine the Unigene ORF, and then determining its CDS and the corresponding  
55 amino acid sequences according to the codon table. At the same time, **TransDecoder**  
56 (version 3.0.1) was used to predict CDS sequences of the un-aligned Unigenes. GO  
57 (Gene Ontology database) functional annotation information was obtained according  
58 to transcripts annotation results of SwissProt and TrEMBL. **KAAS** (version 2.1)  
59 (KEGG Automatic Annotation Server) was used for KEGG (Kyoto Encyclopedia of  
60 Genes and Genomes) annotation.

61

#### 62 1.5 RNA-seq assessment

63 **Bowtie2** (version 2.3.2) was used for aligning the quality control sequences to  
64 the assembled transcripts, and **RSeQC** (version 2.6.1) was used for statistics the  
65 aligned result. Then **BEDTools** (version 2.26.0) was used for homogeneity  
66 distribution check and statistics the gene coverage ratio, **RSeQC** (version 2.6.1)  
67 software was used for duplicate reads analysis.

68

#### 69 1.6 Expression quantity analysis

70 The direct expression of a gene expression level is the abundance of its transcript;  
71 higher transcript abundance represents higher gene expression level. Transcripts Per  
72 Million (TPM) is a measure to calculate the proportion of a transcript in the RNA pool.  
73 It takes into account the sequence depth and the length of the gene as well as the  
74 influence of the sample on the reads count. **Salmon** (version 0.8.2) was used to  
75 calculate the reads count and TPM of unigenes. Differentially expressed genes were  
76 calculated based on the reads count of each gene.

77 For the samples without biological repetition, TMM was used to standardize  
78 the read count data, and then **DEGseq** (version 1.26.0) was used for differently  
79 analysis. For the samples with biological repetition, **DESeq** (version 1.12.4) was  
80 used for analysis. In order to obtain the significant differential genes, the screening  
81 conditions were set as follows: q-value <0.001 and difference multiple  
82 |FoldChange| >2.

83

#### 84 1.7 Relationship analysis of samples

85 The gene expression correlation between samples is an important index to test  
86 whether the experiment is reliable and the sample selection is reasonable. The  
87 closer the correlation coefficient gets to 1, the more similarity the expression  
88 patterns between samples.

89 Principal Component Analysis (PCA) could reflect the distance and difference  
90 between samples through different sample species and function composition  
91 analysis, and the related graph was constructed by **R vegan** package.

92 Principal co-ordinates analysis (PCoA) is a visualization method to study the  
93 similarity or difference of data, and the difference between individuals or groups  
94 can be observed by PCoA. The related graph was constructed by **R vegan** package.

95 Non-metric multidimensional scale (NMDS) is used to map, analyze and  
96 classify the research objects (samples or variables) of the multidimensional space  
97 into the low-dimensional space, while preserving the original relationship between  
98 the objects. The related graph was performed by **R vegan** package.

#### 100 1.8 Functional enrichment analysis

101 **TopGO** (version 2.24.0) was used for GO (Gene Ontology) enrichment, and the  
102 function was thought to be a significant enrichment when the correct p-value (q-value)  
103  $< 0.05$ . The basic unit of GO is GO-term. GO enrichment analysis provides all GO  
104 terms that significantly enriched in DEGs comparing to the genome background.

105 **ClusterProfiler** (version 3.0.5) was used for Kyoto Encyclopedia of Genes and  
106 Genomes (KEGG) enrichment analysis. Pathway enrichment analysis identified  
107 significantly enriched metabolic pathways or signal transduction pathways in DEGs  
108 comparing with the whole genome background.

109 Protein-protein interaction network was performed with **R igraph** package.  
110 STRING is a protein interaction database developed by EMBL, which has collected  
111 the most powerful experimental verification, data mining and homologous  
112 prediction of protein interactions. The differentially expressed genes were mapped  
113 to the responding protein-protein interaction network to extract sub network and  
114 make it visualization, screening key genes according to the topology of genes in the  
115 whole network.