

Comprehensive identification of the full-length transcripts and alternative splicing related to the secondary metabolism pathways in the tea plant (*Camellia sinensis*)

Dahe Qiao¹, Chun Yang¹, Juan Chen¹, Yan Guo¹, Yan Li¹, Suzhen Niu², Kemei Cao² & Zhengwu Chen^{1*}

¹ Tea Research Institute, Guizhou Academy of Agricultural Science, Guiyang 550006 Guizhou, China. ² College of tea science, Guizhou University, Guiyang 550025 Guizhou, China

*Correspondence: Z.W.C, Email: zwchentea@163.com, Tel.: +86- 0851-83761972

Additional Information

Figure. S1 Flow diagram of the experimental design and joint analysis for PacBio sequencing and Illumina sequencing.

Figure. S2 The results of transcriptome integrity assessment based on BUSCO using an embryophyta gene set. The total number of embryophyta gene sets used in this evaluation was 1,440.

Figure. S3 Differential expression analysis of transcripts. (A) Correlation analysis of the expression between two samples. (B) The MA plot of the differentially expressed transcripts. Each point represents a transcript. The x-axis represents A value, which is the logarithm of the mean FPKM values of the two samples, and the y-axis represents M value, which is the logarithm of fold change (FC) in gene expression between two samples. The green dots and red dots represent down and up-regulated differentially expressed transcripts, respectively, and the black dots represent non-differentially expressed transcripts.

Figure. S4 Functional distribution of COG annotation of the transcripts.

Figure. S5 Correlation analysis of gene expression and secondary metabolite accumulation. The secondary metabolism-related structural genes with AS transcripts identified in this study were selected and the correlation analysis were performed on using the Gene Expression and Metabolite accumulation Correlation Analysis Tool (<http://tpia.teaplant.org/Gene2Metabolite.html>).

Figure. S6 Sequence alignment of the amplified fragments of the AS isoforms. The primer sequences were underlined in red.

Table S1. Comparison of single molecule sequencing data of tea plant between two PacBio sequencing platforms.

Table S2. The differential expression transcripts identified in this study.

Table S3. The character of AS events identified in this study.

Table S4. The transcripts involved in flavonoid biosynthesis.

Table S5. The AS events of the transcripts involved in flavonoid, theanine and caffeine biosynthesis.

Table S6. The transcripts involved in theanine biosynthesis.

Table S7. The transcripts involved in caffeine biosynthesis.

Table S8. The character of fusion transcripts identified in this study.

Table S9. Correlation analysis of the gene expression and secondary metabolite accumulation.

Table S10. The primers used in this study.

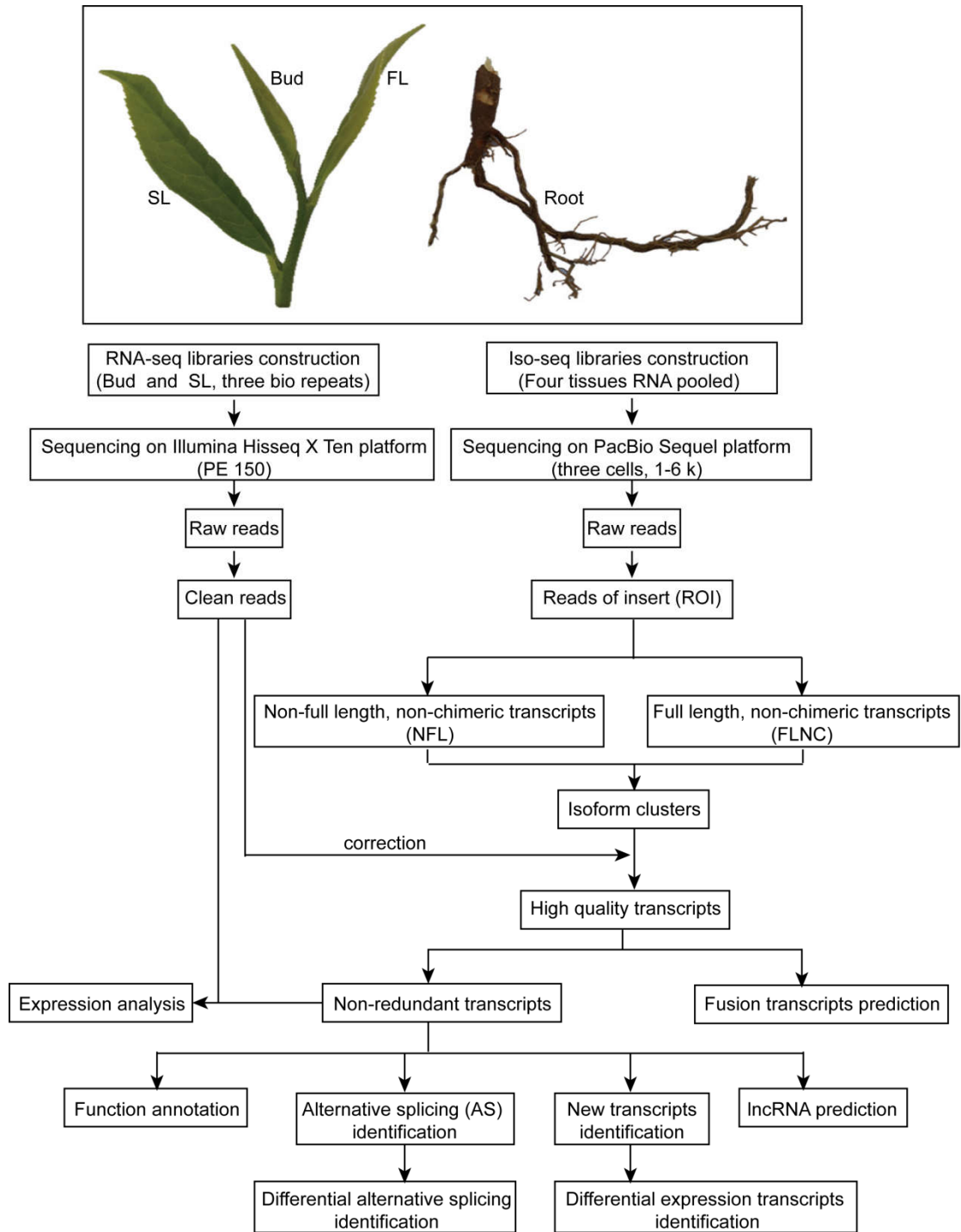


Figure S1. Flow diagram of the experimental design and joint analysis for PacBio sequencing and Illumina sequencing.

BUSCO Assessment Results

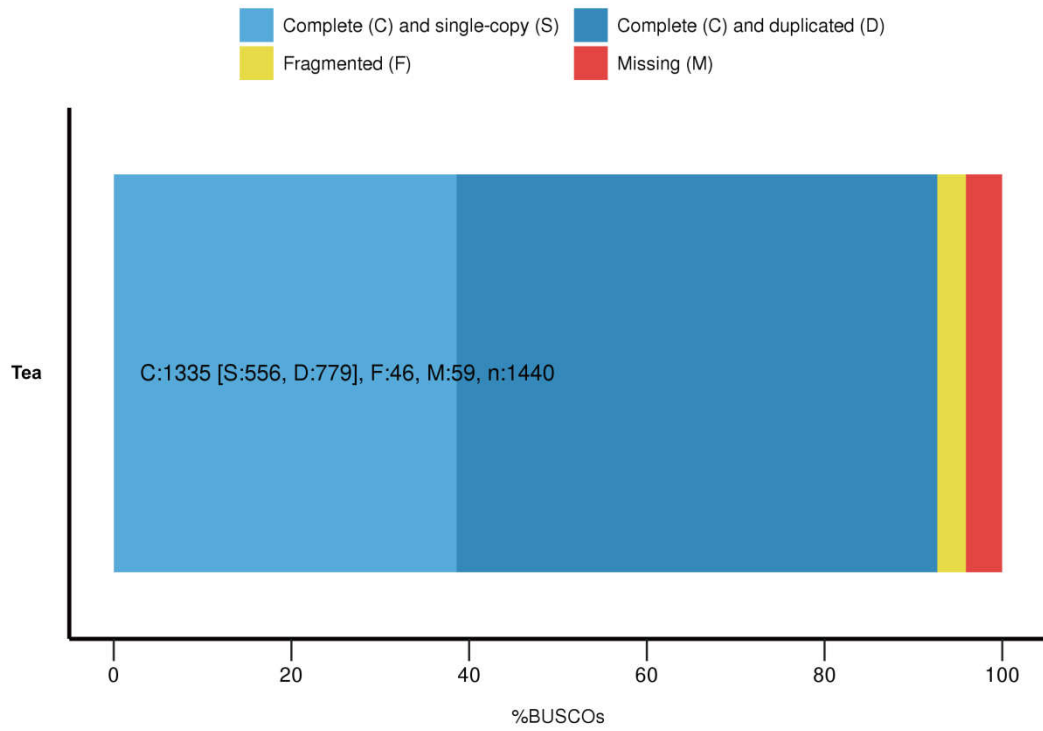


Figure S2. The results of transcriptome integrity assessment based on BUSCO using an embryophyta gene set. The total number of embryophyta gene sets used in this evaluation was 1,440.

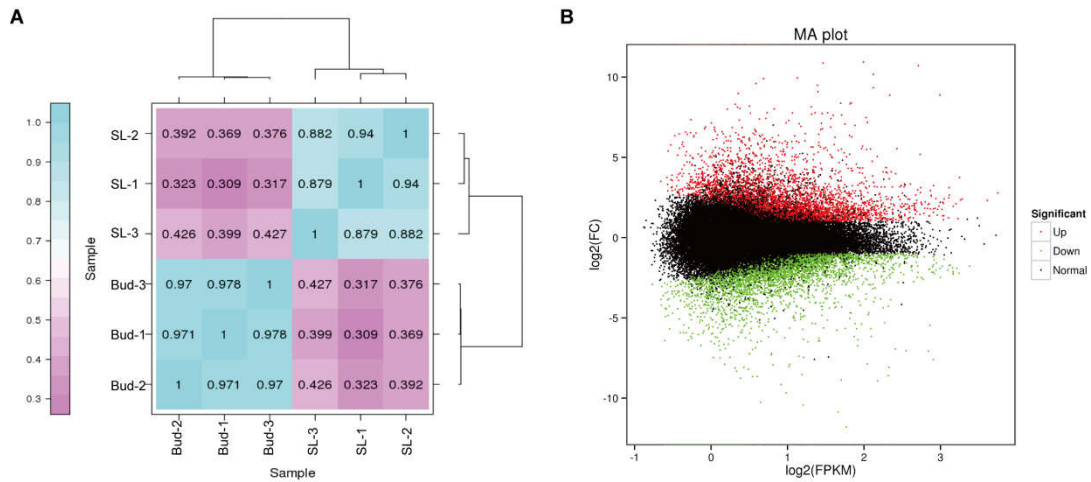


Figure S3. Differential expression analysis of transcripts. (A) Correlation analysis of the expression between two samples. (B) The MA plot of the differentially expressed transcripts. Each point represents a transcript. The x-axis represents A value, which is the logarithm of the mean FPKM values of the two samples, and the y-axis represents M value, which is the logarithm of fold change (FC) in gene expression between two samples. The green dots and red dots represent down and up-regulated differentially expressed transcripts, respectively, and the black dots represent non-differentially expressed transcripts.

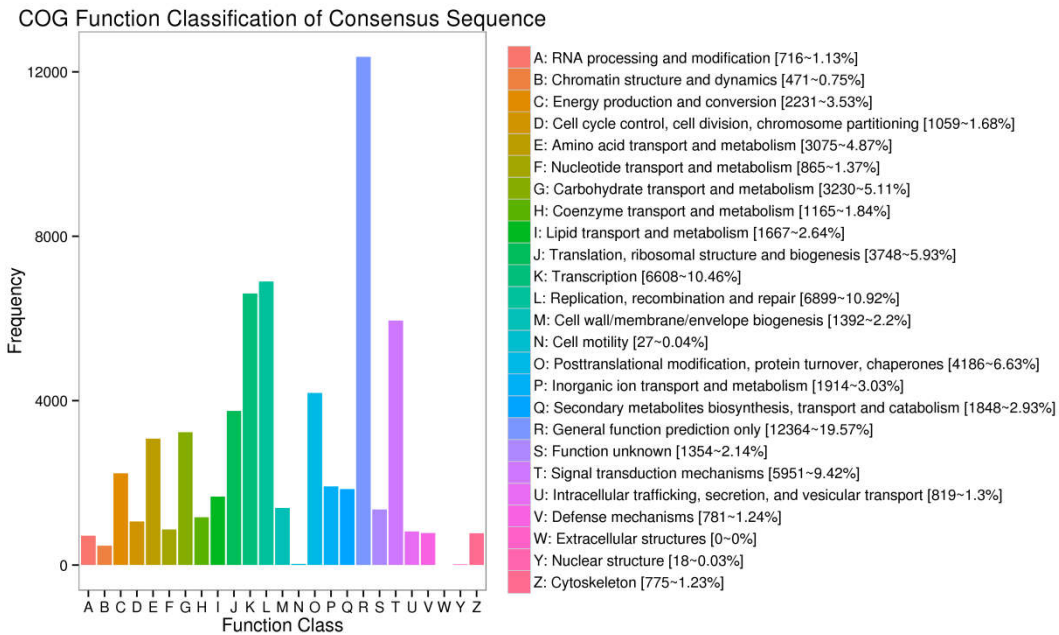


Figure S4. Functional distribution of COG annotation of the transcripts.

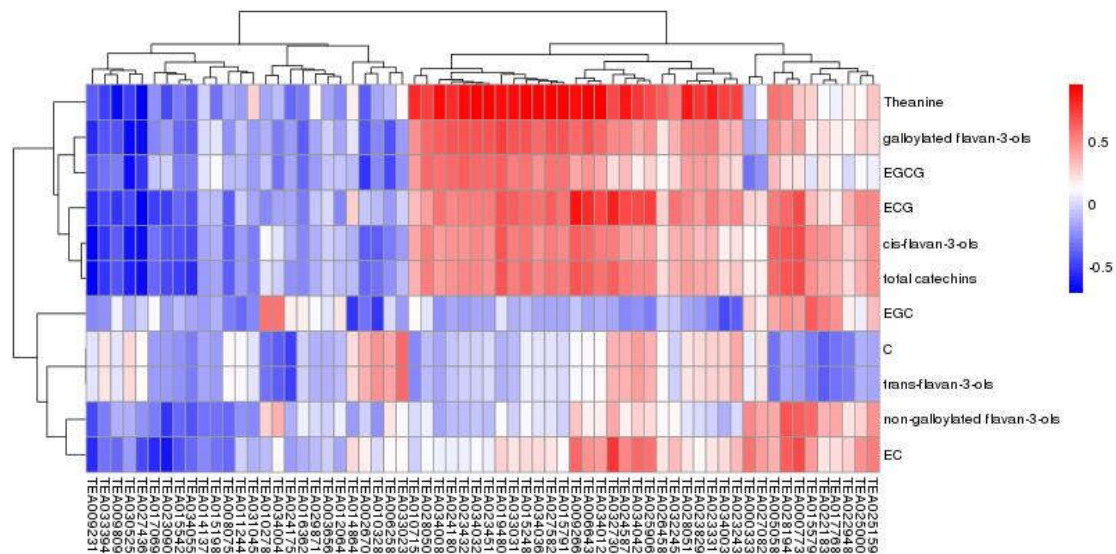


Figure S5. Correlation analysis of gene expression and secondary metabolite accumulation. The secondary metabolism-related structural genes with AS transcripts identified in this study were selected and the correlation analysis were performed on using the Gene Expression and Metabolite accumulation Correlation Analysis Tool (<http://tpia.teaplant.org/Gene2Metabolite.html>).

Figure S6. Sequence alignment of the amplified fragments of the AS isoforms. The primer sequences were underlined in red.

(A) Seq_0:PB.1604.5, Seq_1:PB.1604.2

Seq_0 ----- 0
Seq_1 GGTGTGGAGGTAGTAGACGCACACGCTCACAATCTGGTGGCTCTGGACTCCACTCTTCTTCTCCAATGCTTTTCTGA 80

Seq_0 ----- 0
Seq_1 AGCCTACGGCGATGCTTTATTGCTCGCACCCACGCTCTCAACTTCAAGAGAGGTATAAGGGATATTGCTGAACTGTATG 160

Seq_0 ----- 0
Seq_1 GATCTGAGTTATCCTTGGATGGCATTCAAAAATACCGCAAGGGCAATGGATTGCAATCCATAAGCTCAATATGCTTCAAG 240

Seq_0 ----- 0
Seq_1 GCTGCAAGAATCGCTGCAATACTCATTGATGATGGAATTGAGTTTGACAAAATGCATGACATTGAATGGCATAGGAATTT 320

Seq_0 ----- 0
Seq_1 TGCACCGGTGGTTGGTAGAATATTGAGAATTGAGCATCTGGCTGAGAAGATTCTTGATGAAGGGAGGCCAGATGGATCTA 400

Seq_0 ----- TCACTGAAACATTTATCGGAAAGTTGAAGTCAGTTGCTAATAAAATTGTTGGCTTGAAAAGC 62
Seq_1 CCTGGACATTGGACAGTTTCACTGAAACATTTATCGGAAAGTTGAAGTCAGTTGCTAATAAAATTGTTGGCTTGAAAAGC 480

Seq_0 ATAGCTGCATACCGCAGTGGTCTTGAGATTAATACAAATGTCAACAAGGAAGGAGGCTCAAGCGGTCTTGTGGAAGTTTT 142
Seq_1 ATAGCTGCATACCGCAGTGGTCTTGAGATTAATACAAATGTCAACAAGGAAGGAGGCTCAAGCGGTCTTGTGGAAGTTTT 560

Seq_0 AAATGCTGGGAGCCCCGTTTCGTATCACAAATAAAAACCTCATTGACTATCTCTTCGTGCAGAGTTTGGAGGTTGCCATAC 222
Seq_1 AAATGCTGGGAGCCCCGTTTCGTATCACAAATAAAAACCTCATTGACTATCTCTTCGTGCAGAGTTTGGAGGTTGCCATAC 640

Seq_0 AATATGATTTGCCAATGCAGATACACACTGGTTTTGGAGATAAAGATTTGGATTTAAGGCTCTCCAATCCCTGCATCTC 302
Seq_1 AATATGATTTGCCAATGCAGATACACACTGGTTTTGGAGATAAAGATTTGGATTTAAGGCTCTCCAATCCCTGCATCTC 720

Seq_0 CGCACCTTCTTGAGGACAAGAGATTCTCTAAGTGCCGCTTAGTACTTTTACATGCATCATACCCATTTTCAAAGGAAGC 382
Seq_1 CGCACCTTCTTGAGGACAAGAGATTCTCTAAGTGCCGCTTAGTACTTTTACATGCATCATACCCATTTTCAAAGGAAGC 800

Seq_0 ATCATATCTAGCCTCCATTTATTTCTCAGGTTTACCTTGATTTGGTTGGCTGTTC 438
Seq_1 ATCATATCTAGCCTCCATTTATTTCTCAGGTTTACCTTCATTTGGTTGGCTGTTC 856

(F) Seq_0:PB.6810.2, Seq_1:PB.6810.3

Seq_0 **TTCCGATCTTTGCAACTTGA**ATCTATCTGAGTC**ACCGAGAAGATCATTGCAGAGTACATATGGATCGGTGGATCTGGTA** 80
Seq_1 **TTCCGATCTTTGCAACTTGAATCTATCTGAGTC****ACCGAGAAGATCATTGCAGAGTACATATGGATCGGTGGATCTGGTA** 80

Seq_0 **TGGACCTCAGA**AAGCAAAGCCAGGACCTGAATGCACCAGTCTCTGATCCTTCAAAGTTACCACAATGGAACCTACGATGGT 160
Seq_1 **TGGACCTCAGA**----- 91

Seq_0 TCCAGCACTGGCCAAGCCCCTGGCGAGGACAGTGAAGTGATTTTATATCCCAGGCAATTTATAAGGACCCATTTCAGGAG 240
Seq_1 ----- 91

Seq_0 **AGGCAACAACATTTCTTGTAAATGTGTGATGCTTACACGCCGGGTGGAGAGCCAATCCCAACAATAAGAGGTTTGATGCTG** 320
Seq_1 **-GGCAACAACATTTCTTGTAAATGTGTGATGCTTACACGCCGGGTGGAGAGCCAATCCCAACAATAAGAGGTTTGATGCTG** 170

Seq_0 **CCAAGATATTCAGCCACCCTGATGTTGCTGAGGAACCTTGGTATGGTATAGAGCAGGAGTACACTTTGTTGCAGAAA** 400
Seq_1 **CCAAGATATTCAGCCACCCTGATGTTGCTGAGGAACCTTGGTATGGTATAGAGCAGGAGTACACTTTGTTGCAGAAA** 250

Seq_0 **GAAGTGAAGTGGCCGATTGGTTGGCCTGTGGGAGGTTATCCTGGACCACAGGGACCATACTACTGTGGTATTGGTCCGGA** 480
Seq_1 **GAAGTGAAGTGGCCGATTGGTTGGCCTGTGGGAGGTTATCCTGGACCACAGGGACCATACTACTGTGGTATTGGTCCGGA** 330

Seq_0 **TAAAGCTTTTGGGCGAGACATTGTCGATGCCATTATAAAGCATGTCTTTATGCTGGTATTAACATTAGTGGCATCAATG** 560
Seq_1 **TAAAGCTTTTGGGCGAGACATTGTCGATGCCATTATAAAGCATGTCTTTATGCTGGTATTAACATTAGTGGCATCAATG** 410

Seq_0 **GAGAGGTGATGCCGGGTCAGTGGGAATCCAAGTTGGGCCTTCTGTTGGCATCAGTTCTGGAGATCAGTTGTGGATGGC** 639
Seq_1 **GAGAGGTGATGCCGGGTCAGTGGGAATCCAAGTTGGGCCTTCTGTTGGCATCAGTTCT****GGAGATCAGTTCTGGATGGC** 489

Table S1. Comparison of single molecule sequencing data of tea plant between two PacBio sequencing platforms.

Platform	Library Size	Cell Number	Reads of Insert	Read Bases of Insert (bp)	Mean Read Length of Insert (bp)	Mean Read Quality of Insert	Mean Number of Passes
PacBio Sequel (This study)	1-6K	3	1,388,066	2,446,983,161	1,762	0.92	10
	<1K	2	55,037	42,298,256	768	0.96	33
PacBio RS II (Xu et al. 2017)	1-2K	2	135,732	293,235,156	2,160	0.94	11
	2-3K	2	119,629	361,701,882	3,023	0.91	9
	3-6K	1	51,549	200,280,602	3,885	0.88	4