1 ## *De novo* transcriptome sequence of *Senna tora* provides insights

2 ## into anthraquinone biosynthesis

3

4 Sang-Ho Kang[1,*, †], Woo-Haeng Lee[2,*], Chang-Muk Lee[3], Joon-Soo Sim[3], So Youn Won[1],

5 So-Ra Han[2], Soo-Jin Kwon[1], Jung Sun Kim[1], Chang-Kug Kim[1,†], Tae-Jin Oh[2,4,5,†]

6

7 [1]Genomics Division, National Institute of Agricultural Sciences, RDA, Jeonju, Korea; [2]Department of

8 Life Science and Biochemical Engineering, SunMoon University, Asan, Korea; [3]Metabolic

9 Engineering Division, National Institute of Agricultural Sciences, RDA, Jeonju, Korea; [4]Genome-

10 based BioIT Convergence Institute, Asan, Korea; [4]Department of Pharmaceutical Engineering and

11 Biotechnology, SunMoon University, Asan, Korea

12

13 * These authors contributed equally to this work.

14

15 [†]Corresponding author: Sang-Ho Kang; Chang-Kug Kim; Tae-Jin Oh

16 E-mail: hosang93@korea.kr; Tel: +82-63-238-4570; Fax: +82-63-238-4554

17 E-mail: chang@korea.kr; Tel: +82-63-238-4555; Fax: +82-63-238-4554

18 E-mail: tjoh3782@sunmoon.ac.kr; Tel: +82-41-530-2677; Fax: +82-41-530-2279

19

# **Abstract**

*Senna tora* is an annual herb with rich source of anthraquinones that have tremendous pharmacological properties. However, there is little mention of genetic information for this species, especially regarding the biosynthetic pathways of anthraquinones. To understand the key genes and regulatory mechanism of anthraquinone biosynthesis pathways, we performed spatial and temporal transcriptome sequencing of *S. tora* using short RNA sequencing (RNA-Seq) and long-read isoform sequencing (Iso-Seq) technologies, and generated two unigene sets composed of 118,635 and 39,364, respectively. A comprehensive functional annotation and classification with multiple public databases identified array of genes involved in major secondary metabolite biosynthesis pathways and important transcription factor (TF) families (MYB, MYB-related, AP2/ERF, C2C2-YABBY, and bHLH). Differential expression analysis indicated that the expression level of genes involved in anthraquinone biosynthetic pathway regulates differently depending on the degree of tissues and seeds development. Furthermore, we identified that the amount of anthraquinone compounds were greater in late seeds than early ones. In conclusion, these results provide a rich resource for understanding the anthraquinone metabolism in *S. tora*.

**Key words**

*Senna tora*, anthraquinone, secondary metabolite, transcriptome analysis, transcription factor

# Introduction

41       *Senna tora* (Subfamily, Caesalpiniaceae; and Family, Leguminosae) also known as

42    *Cassia tora*, is an annual xerophytic shrub which grows in the arid zones after the rainy

43    season [1]. This plant is mostly found in India, China, Sri Lanka, Nepal, the Korean

44    peninsula, and other Asian countries. Its name varies in different locales such as Foetid Senna

45    tora, Sickle senna, Wild senna, Coffee pod, Tovara, Chakvad, and Ringworm plant. *S. tora*

46    leaves, seeds, and roots have long been used as food ingredients. It is also valued as a

47    medicinal plant in Ayurveda, commonly used as a depurative, antiperiodic, anthelmintic,

48    liver tonic, hepatic disorders, dyspepsia leprosy, constipation, intermittent fever, cough,

49    bronchitis, ringworm infection, ophthalmic, skin diseases, and others [2, 3]. It has also been

50    used as laxative and a tonic, and is popularly served as a roasted tea throughout Korea and

51    China [4]. The seeds of *S. tora* contain a variety of bioactive anthraquinone substances,

52    including chrysophanol, obtusin, obtusifolin, aurantio-obtusin, chyro-obtusin, obstsifolin,

53    emodin, rubrofusarin, gentibioside, and rhein. Chryophanol is primarily responsible for the

54    plant's pharmacological properties [5, 6]. *S. tora* mainly contains anthraquinone glycosides

55    and flavonoids [7]. Recently, *S. tora* seed extract (STE) and its active compound aurantio-

56    obtusin have been found to suppress degranulation, histamine production, and reactive

57    oxygen species generation, and also to inhibit the production and mRNA expression of

58    cyclooxygenase 2. STE and aurantio-obtusin also suppressed IgE-mediated FcεRI such as

59    phosphorylation of Syk, protein kinase Cμ, phospholipase Cγ, and extracellular signal-

60    regulated kinases. This suggests that STE and aurantio-obtusin can be beneficial to the

61    treatment of allergy-related diseases [8].

62       Anthraquinones, secondary metabolites occurring in bacteria, fungi, lichens, and

63    higher plants, seem to originate from a variety of different precursors and pathways. There

3

64 are two pathways leading to anthraquinone biosynthesis in higher plants: the polyketide

65 pathway and the chorismate/*O*-succinylbenzoic acid pathway. The latter occurs in the plant

66 family Rubiaceae and synthesizes aromatic compounds known for a broad spectrum of

67 bioactivity, such as anticancer, cathartic, anti-inflammatory, anti-microbial, diuretic,

68 vasorelaxing, and phytoestrogen activities, and has recently shown therapeutic potential in

69 autoimmune diabetes [9]. Emodin, physicion, aloe-emodin, and rhein isolated from *S. tora*

70 seed shows antifungal properties against phytopathogenic fungi [10]. Likewise, rhein shows

71 high antibacterial activity towards *Porphyromonas gingivalis* and synergistic antibacterial

72 activity with metronidazole or natural compounds, and the recent studies suggest the

73 immunomodulatory activity of rhein [11-13]. The extract of *S. tora* is found to have

74 hypolipidemic activity, hepatoprotective, and antioxidant effects [2, 14, 15]. Anthraquinones

75 from *S. tora* exhibit significant inhibitory properties against angiotensin-converting enzyme

76 (ACE). Among the various anthraquinones, only anthraquinone glycoside demonstrates

77 marked inhibitory activity against ACE [16].

78      RNA sequencing (RNA-Seq), a technology that can be used to profile the complete

79 gene space of various organisms due to their high throughput, accuracy, and reproducibility,

80 has accelerated the discovery of new genes or analysis of tissue-specific and functional

81 expression patterns in large, complex genomes like those of plants [17-19]. But in the

82 absence of reference genome information considerable small transcripts hinder the accuracy

83 of the construction of RNA sequencing libraries and the efficiency of functional gene

84 prediction or annotation. Short-length RNA sequencing data limit the creation of a longer,

85 accurate contig assembly, resulting in chimeric contigs and/or low gene annotation [20].

86 Moreover, small laboratories require high sequencing costs due to the need for long reads and

87 high-depth short read sequences to be accurate in *de novo* assembly. Plants with large

4

88　genomes pose even more difficult as in, for example, the common soybean crop, which has a

89　genome size of ~1.1Gb [21]. To improve the comprehensive accuracy of gene prediction,

90　there is a need to introduce a new approach, the "Isoform sequencing (Iso-Seq)." Thanks to

91　its long-read technology, Iso-Seq facilitates identifying new isoforms with a high level of

92　accuracy [22]. Advances in technology enable long reads in the range of 1.5-10 kb, which are

93　able to provide full-length mRNA isoforms, detect new isoforms, and skip the transcript

94　reconstruction process by identifying isoforms directly [23]. In this study, we present the

95　transcriptome analysis of the plant *S. tora* from 4 different sources using RNA-Seq and Iso-

96　Seq, providing insights of key genes involved in anthraquinone biosynthesis in the

97　pharmacologically important herb *S. tora*.

98

# Materials and methods

## Plant material and RNA preparation

101　　　Specimens of *S. tora* (cv. Myeongyun) grown in an experimental plot of National

102　Institute of Horticultural and Herbal Science (Eumseong) field were used for transcriptome

103　analysis. Leaf, root, and early- and late-stage seed tissues were harvested from healthy plants,

104　and stored at -80°C until used for RNA extraction. Total RNA was extracted from leafs,

105　roots, and two stages of seeds of *S. tora* using the RNeasy Plant Mini kit (Qiagen, InS.,

106　Valencia, CA, USA). RNA purity was determined using NanoDrop8000 Spectrophotometer

107　and Agilent Technologies 2100 Bioanalyzer, and total RNA integrity was identified as having

108　a minimum integrity value of 7.

109

## Illumina short-read sequencing

111       The poly (A)$^+$ mRNA was purified and fragmented from 1 µg of total RNA using

112   poly-T oligo-attached magnetic beads by two rounds of purification. Using reverse

113   transcriptase, random hexamer primers, and dUTP, the randomly-cleaved RNA fragments

114   were transcribed reversely into first-strand cDNA. A single A-base was added to these cDNA

115   fragments followed by adapter ligation. The products were purified and concentrated by PCR

116   in order to generate a final-strand specific cDNA library. The quality of the amplified

117   libraries was verified using capillary electrophoresis (Bioanalyzer, Agilent). Quantitative

118   PCR (qPCR) was carried out using SYBR Green PCR Master Mix (Applied Biosystems).

119   Then we pooled together equimolar amounts of libraries that were index-tagged. The cBot-

120   automated cluster creation system (Illumina) performed cluster generation in the flow cell.

121   The sequencing was performed with 2 x 100 bp read length of the flow cell loaded on a

122   HiSeq 2500 sequencing system (Illumina).

123

## Long-read sequencing

125       Libraries for Pacific Biosciences Single Molecule Real Time (SMRT) sequencing

126   were prepared from the aforementioned cDNAs. Cycle optimization was performed to

127   determine the optimal number of cycles for large-scale PCR. We prepared 3 fraction cDNAs

128   (1-2 kb, 2-3 kb, and 3-6 kb) using the BluePippin Size selection system. The SMRTbell

129   library was constructed by using SMRTbell$^{TM}$ Template Prep Kit (PN 100-259-100). The

130   DNA/Polymerase Binding Kit P6 (PacBio) was used for DNA synthesis after the sequencing

131   primer annealed to the SMRTbell template. Following the polymerase binding reaction, the

132   MagBead Kit was used to bind the library complex with MagBeads before sequencing.

133   MagBead-bound cDNA complexes result in increased number of reads per SMRT cell. This

134   polymerase-SMRTbell-adaptor complex was then loaded into zero-mode waveguides

135   (ZMWs). The SMRTbell library was sequenced using 8 SMRT cells (Pacific Biosciences)

136   with C4 chemistry (DNA sequencing Reagent 4.0). 1 × 240 minute movies were captured for

137   each SMRT cell using the PacBio RS II sequencing platform.

138

## De novo transcriptome assembly and sequence clustering

140       Raw data of the *S. tora* transcriptome generated from Illumina HiSeq were

141   preprocessed to remove nonsense sequences including adaptors, primers, and low quality

142   sequences (Phred quality score of less than 20) using NGS QC Toolkit [24]. The raw data

143   were further processed to remove ribosomal RNA using riboPicker v0.4.3 [25]. The

144   preprocessed reads were then assembled using Trinity [26]. Assembly statistics were

145   calculated using in-house Perl scripts. Assembled transcripts were clustered (CD-HIT-EST

146   v4.6.1) [27] in order to reduce sequence redundancy. Sequence identity threshold and

147   alignment coverage (for the shorter sequence) were both set as 90% to generate clusters. Such

148   clustered transcripts are defined as reference transcripts in this work.

149

## Illumina expression quantification and differential expression analysis

152       The cleaned reads from each tissue were aligned with the abundant transcriptome

153   assembly using Bowtie2 [28]. The aligned reads were quantified as fragments per million

154   reads (FPKMs) against non-redundant combined transcript sequences (at 90% sequence

155   similarity by CD-HIT-EST). The reads counting of alignments was performed using RSEM

156   (RNA-Seq by Expectation Maximization)-1.2.25 [29]. The differential expression analysis

157   was performed using the DESeq2 packages [30]. Differentially expressed genes (DEGs) were

158    identified using the combined criteria of a more than twofold change and significance with P-

159    value threshold of 0.001 based on the three biological replicates.

160

## Functional annotation and classification

162    All the assembled unigenes were annotated by BLAST program [31] against the

163    National Center for Biotechnology Information (NCBI) nonredundant (Nr) protein database,

164    the Swiss-Prot protein database, and the Kyoto Encyclopedia of Genes and Genomes

165    (KEGG) pathways database with an E-value cutoff of $10^{-5}$. The best aligning results were

166    selected to annotate the unigenes. Whenever the aligning results from different databases

167    conflicted, the results from Swiss-Prot database were preferentially selected, followed by Nr

168    database and KEGG database. Functional categorization by Geno Ontology (GO) terms [32]

169    was carried out by Blast2GO program [33] with E-value threshold of $10^{-5}$. AgriGO [34] was

170    used to determine over-representation of GO categories (e.g., biological processes).

171

## Identification of transcription factor families

173    To investigate the putative transcription factor families in *S. tora*, unigenes were

174    mapped against all the transcription factor protein sequences made available by the Plant

175    Transcription Factor Database (PlantTFDB 4.0; http://planttfdb.cbi.pku.edu/download.php)

176    using BLASTX with E-value threshold of $10^{-5}$.

177

## Quantitative RT-PCR analysis

179    Total RNA was extracted by using the RNeasy Plant Mini Kit (Qiagen, Valencia,

180    CA, USA) following the manufacturer's instructions. The quality of the isolated RNA was

181 checked on ethidium bromide-stained agarose gels, and its concentration was calculated

182 according to the measured optical density (OD) of the samples at 260 and 280 nm

183 (DropSense96C Spectrophotometer, Trinean, Belgium). The 1 μg of the total RNA was used

184 for the cDNA synthesis using SuperScript™ III first strand RT-PCR kit (Invitrogen,

185 Carlsbad, CA, USA) with an oligo(dT)$_{20}$ primer. After cDNA was obtained from *S. tora*,

186 qRT-PCR was performed using gene-specific primers (S1 Table). Real-time PCR analysis

187 was optimized and performed using the Roche LightCycler® 480 II instrument and SYBR®

188 Green Real-Time PCR Master Mix (Bio-Rad, InS., Hercules, CA, USA) under condition of

189 an initial denaturation at 95°C for 30 s followed by 40 cycles of denaturation at 95°C for 10 s,

190 annealing and extending at 55°C for 15 s. The relative expression of specific genes was

191 quantified using the $2^{-\Delta\Delta Ct}$ calculation according to the manufacturer's software [35] (where

192 $\Delta\Delta C_t$ is the difference in the threshold cycles), and the internal reference gene was the

193 elongation factor 2 for data normalization. Reliability of the amplification parameters was

194 analyzed at 1:15 dilutions of the cDNA samples. The mean threshold cycle values for the

195 genes of interest were calculated from three experimental replicates.

196

## Extraction of anthraquinones and LC-MS analysis

198 Early- and late-stage of seed samples were extracted with methanol using sonication

199 for 30 min at 60°C. After extraction, samples were centrifuged at 12,000 rpm for 3 min at

200 25°C and the supernatant was filtered with 0.2 μm Acrodisc® MS Syringe Filters with

201 WWPTFE membrane (Pall Corporation, Port Washington, NY, USA). Quantitative analysis

202 of anthraquinones was performed by a Triple TOF 5600+ Spectrometer with a DuoSpray ion

203 source (AB Sciex, Ontario, CA, USA) coupled with a Nexera X2 UHPLC (Shimadzu, Kyoto,

204 Japan) equipped with binary solvent manager, sample manager, column heater, and

205   photodiode array detector. UHPLC was performed on a ACQUITY UPLC®BEH C18

206   column (1.7 μm, 2.1 x 100 mm, Waters Corporation, Milford, USA) and mobile phases

207   consisted of 5 mM ammonium acetate in water (eluent A) and 100% acetonitrile (eluent B).

208   The gradient profile was as follows: 0-1 min, 20% B; 1-3.5 min, 10-30% B; 3.5-8 min, 30-

209   50% B; 8-12 min, 50-100% B; 11-17 min, 100% B. The flow rate was 0.5 mL/min and five

210   microliters of samples were injected. For detecting peaks from test samples, MS parameter in

211   ESI-negative mode was used as follows: nebulizing gas, 50 psi; heating gas, 50 psi; curtain

212   gas, 25 psi; desolvation temperature, 500°C; ion spray voltage floating, 4.5 kV.

213

## Data availability

215      The RNA-Seq and Iso-Seq sequences generated from Illumina and PacBio RS II

216   sequencing of four tissue samples of *S. tora* were deposited at the National Center for

217   Biotechnology Information (NCBI) Sequence Read Archive database with the accession

218   number SRP159435.

219

# Results and discussion

## RNA sequencing and de novo transcriptome assembly

222      *De novo* transcriptome analysis is a good tool for generating the overall genetic

223   information of an organism without full genome sequencing and leads to discoveries of new

224   genes, molecular markers, and tissue-specific expression patterns. We used the Illumina

225   HiSeq 2500 system and PacBio RS II platform to sequence the cDNA libraries of the leaf,

226   root, and early- and late-stages of seed for elucidating secondary metabolites biosynthesis and

227   understanding their spatial and temporal expression pattern in *S. tora*. Illumina Hiseq 2500

228     sequencing platform produced 278,031,495 raw reads and averaged 23,169,291 reads per

229     tissue (S2 Table). In total, more than 270 million reads showed high quality read rates (Q30

230     values) of over 88.00% (S2 Table). The Trinity assembler from the four different libraries

231     generated a total of 118,635 unigenes that were more than 300 base pairs (bp) long (Fig 1).

232     The length of the transcripts varied from 300 to 18,622 bp with an average length of 832.25

233     bp, the N50 length of 1,082 bp, and the GC content of 39.51% (Table 1).

234          A unigene, the assembled transcript that represents a hypothetical gene, can be

235     represented by several isomers as different forms of the same protein. The PacBio RS II

236     sequencing platform produced 768,745 raw reads. After classification and clustering, 118,703

237     high-quality isoforms were obtained from three different libraries, which contained 39,672,

238     32,954, and 46,077 high-quality isoforms per library sizes (<2 kb, 2-3 kb, and >3 kb) (S3

239     Table). The 118,703 high-quality isoforms from three different libraries generated 39,364

240     non-redundant unigenes after the CD-HIT-EST program removed redundant isoforms. The

241     total size of the assembly was 112 MB with 57% of transcripts larger than 500 bp and 12%

242     larger than 2,000 bp. In total, our analysis generated two unigene sets: 118,635 from RNA-

243     Seq and 39,364 from Iso-Seq (Fig 1). The two unigene sets showed similar GC contents.

244     However, overall unigene lengths of each set showed that the length of the Iso-Seq was

245     longer than RNA-Seq. Unigenes obtained from Iso-Seq were better in terms of minimum

246     length, average length, and N50 length (Table 1).

247          In our analyses, we used the Iso-Seq unigene set mainly as a reference for RNA-Seq

248     data. Due to other dissimilar characteristics, such as the transcript length between the RNA-

249     Seq and Iso-Seq gene sets, this study did not constitute an integrated unigene set. Later, we

250     plan to create one using the reference-guided method when the *S. tora* genome sequencing is

251     completed.

11

252

## Functional annotation and classification

254    Annotation of function is required to characterize transcripts and understand the

255    complexity and diversity of an organism. For the functional annotation, the assembled

256    118,635 unigenes obtained from RNA-Seq of leaf, root, early seed, and late seed tissue

257    samples were screened using an FPKM criterion of $\geq 1$, which resulted in 56,707 unigenes.

258    To obtain the best annotations, assembled 56,707 RNA-Seq unigene sets and 39,364 Iso-Seq

259    unigene sets of *S. tora* were aligned with four public protein databases. We used the

260    BLASTX program against NCBI Nr, Swiss-Prot, KEGG, and GO protein databases with an

261    E-value threshold of 1e-5. Annotations of RNA-Seq and Iso-Seq unigenes resulted in the

262    identification of 43,286 and 36,882 unigene sets, which were respectively matched with

263    known proteins. The Venn diagram displays the unique best BLASTX hits from NCBI Nr,

264    Swiss-Prot, KEGG, and GO databases (S1 Fig). The overlapping regions of the four circles

265    indicate the number of unigenes sharing BLASTX similarities in respective databases. The

266    Venn diagram of RNA-Seq showed significant matches: 32,469 to Swiss-Prot (75.01%),

267    42,552 to NCBI Nr (98.30%), 3,279 to KEGG (7.58%), and 30,287 to GO terms (69.97%).

268    So did the Venn diagram of Iso-Seq: 30,626 to Swiss-Prot (83.04%), 36,830 to NCBI Nr

269    (99.86%), 6,441 to KEGG (17.46%), and 26,762 to GO terms (72.56%). In summary, 43,286

270    RNA-Seq and 36,882 Iso-Seq unigene sets had at least one significant protein match to these

271    databases. The pattern of annotation of RNA-Seq and Iso-Seq showed that the Iso-Seq is

272    better than RNA-Seq at annotating essential data. Non-significant genes that may represent

273    new genes, non-coding RNA, or RNA representing unnecessary information is not evaluated

274    in this annotaion, and further analysis is required. Matches to the Nr database also indicated

275    that a large number of the *S. tora* unigenes closely matched the sequences of *Glycine max*

12

276 (26.94%), *Glycine soja* (13.07%), *Vigna radiate* var. *radiata* (3.21%), *Cicer arietinum*

277 (9.38%), and *Phaseolus vulgaris* (5.63%). Unigenes of 15 species in the Nr database had >

278 1% match with those of *S. tora* (S2 Fig).

279       To further functionally characterize the *S. tora* transcriptome, we classified the functions

280 of RNA-Seq and Iso-Seq unigenes using GO analysis. The distribution of RNA-Seq and Iso-

281 Seq unigene sets in different GO categories is shown in Fig 2. The three main categories of

282 GO annotations of RNA-Seq included 26,616 GO terms (42.12%) for biological process,

283 20,211 terms (31.98%) for molecular function, and 16,365 terms (25.90%) for cellular

284 component. Among biological process, organic substance metabolic process (17.00%) and

285 primary metabolic process (16.00%) were the most abundant GO categories. Regarding

286 molecular function, GO terms related to organic cyclic compound binding (19.00%) and

287 heterocylic compound binding (19.00%) were the most abundant, while cell part (22.00%)

288 and cell (22.00%) were the mostly represented GO categories in cellular components.

289 Conversely, the three main categories of GO annotation of Iso-Seq include 57,137 GO terms

290 (45.64%) for biological process, 31,562 terms (25.13%) for molecular function, and 36,876

291 terms (29.37%) for cellular component. Among biological process, organic substance

292 metabolic process (16.00%) and primary metabolic process (16.00%) were the most abundant

293 GO categories of biological process. The GO terms related to nucleotide binding (16.00%)

294 and nucleoside phosphate binding (16.00%) were the most abundant in molecular function

295 categories. Also, the most abundant GO categories in cellular component were cell part

296 (24.00%) and cell (24.00%). GO terms pattern of RNA-Seq and Iso-Seq was similar in

297 patterns.

298       Transcription factor (TF) families, including ARF, bHLH, bZIP, C2H2, ERF, MIKC,

299 MYB, NAC, and WRKY, play a key regulatory role in the expression of genes, which are

13

300 involved in plant secondary metabolism and response to environmental stress, by binding to

301 specific cis-regulatory elements of the promoter regions. The number of genes encoding for

302 different TF families varies in different plants to perform species-/tissue-specific or

303 developmental stage-specific function [36]. In our study, 3,284 RNA-Seq and 3,576 Iso-Seq

304 were generated with a total of 6,860 unigenes assigned to 56 TF families. Among these,

305 bHLH (521, 15.86%) were found to be the most abundant in RNA-Seq followed by WRKY

306 (243, 7.40%), C2H2 (189, 5.76%), MYB (177, 5.39%), bZTP (170, 5.18%), and NAC (150,

307 4.57%). Similarly, in the Iso-Seq, bHLH were found to be the most abundant followed by

308 WRKY, but the other TF families showed a slight ranking change (S3A Fig).

309 Expression of the gene varies depending on the environment in which each species is

310 exposed, and specific or large amounts of the gene are expressed. The degree of expression of

311 the TF family, which mediates and controls their expression, is essential for the molecular

312 genetics of organisms, so in order to investigate tissue specific gene expression in *S. tora* we

313 studied the expression of genes in leaf, root and early and late seed tissues. Interestingly,

314 different expression patterns for TFs were observed in four tissues of *S. tora*. Some TFs were

315 unique to each tissue, whereas others were enriched in respective tissues. The 35 and 98 TFs

316 among a total of 133 TFs expressed in leaf, 41 and 97 from 138 TFs in root, 30 and 51 from

317 81 TFs in early-stage, and 15 and 18 from 33 TFs in late-stage during seed development were

318 tissue-enriched and -specific (S3B Fig). Notably, growth regulating factor (GRF) in the TF

319 family was dominantly expressed in late-stage seed tissue (S4 Table). GRFs are plant-specific

320 transcription factors that were originally identified for their roles in stem and leaf

321 development [37]. However, recent studies highlight its importance in other central

322 developmental processes including root development, flower, and seed formation. Expression

323 of GRFs has also been observed in various rice and maize tissues, suggesting their

14

324     involvement in seed development [38, 39].

325

## Differential gene expression analysis during seed development

327     To compare genes of *S. tora* with differential expression level in late-stage seed

328     development to those in early-stage development, we used the DESeq method. The

329     transcripts with log2 fold change (FC) >1 and p-value < 1e-3 were considered as

330     differentially expressed genes (DEGs). Pair-wise comparison of transcripts between early-

331     and late-stages of seed development resulted in a total of 14,825 DEGs in RNA-Seq. As

332     seeds matured, 4,935 genes were identified as up-regulated and 9,890 genes were down-

333     regulated. These genes belong to diverse functional groups including glycosyl hydrolases,

334     dehydrogenases, transferases, kinases, phosphatases, cytochrome P450, oxygenases, and

335     hormone-responsive proteins. A heat map was constructed to cluster the top 50 DEGs based

336     on the similarity and diversity of expression profiles using normalized FPKM values within a

337     range of 6 to 16 (Fig 3). Specifically, transcripts of various proteins are expressed differently

338     depending on the tissue and stage of seed. In early-stage seeds, the expression of chalcone

339     synthase, peroxidase, and cell wall/vacuolar inhibitor of fructosidase were higher than those

340     of late-stage seeds. In particular, C/VIF releases glucose and fructose in irreversible

341     reactions, which is essential to plant growth, storage compound accumulation, and stress

342     response [40]. Conversely, in late-seed development, late embryogenesis-abundant (LEA)

343     proteins and heat shock proteins (HSPs) appeared to be more abundant than early seeds like

344     the adlay species [41].

345     Previously, the expression of genes in leaves, roots, and early- and late-seed tissues were

346     examined to investigate the tissue-specific gene expression of *S. tora*. During this process the

347     transcripts exhibiting tissue-specific expression were identified and the top 10 transcripts

15

348    were selected (S4 Fig). Real-time PCR analysis was performed in order to accurately identify

349    differential expression of selected transcripts in the data. Expression analysis was carried out

350    from the selected genes belonging to carbohydrate mechanism, the secondary metabolite

351    pathway, and the associated transcription factors (Fig 4). These results were consistent with

352    tissue-specific gene expression data in various tissues. As results, 3 genes were identified in

353    the qRT-PCR of the seeds to be specifically expressed compared to other tissues. Cell

354    wall/vacuolar inhibitor of fructosidase 2(C/VIF2) play important roles in carbohydrate

355    metabolism, stress responses, and sugar signaling. The specific expression of C/VIF2 in early

356    seeds is implicated in several mechanisms of maturation. Cytochrome P450 83B1 genes

357    showed the highest expression levels in leaf, followed by root, late seed, and early seed.

358    Cytochrome P450 83B1 protein is known to be involved in auxin homeostasis and

359    glucosinolate biosynthesis associated with plant growth and pathogenic responds [42]. Also,

360    seed biotin-containing protein gene showed the highest expression levels in late seed,

361    followed by early seed, demonstrating that the protein plays an important role in the

362    developmental stage of the seed. And organic-cation/carnitine transporter 1 protein gene

363    expressed high levels in root, followed by leaf and late seed. Organic-cation/carnitine

364    transporter families are generally characterized as polyspecific transporters involved in the

365    homeostasis of solutes in animals [43]. Although some publications have suggested that this

366    protein is known as stress-regulated member of plants and that it is involved in plant growth

367    [43], little is known about the function, localization, and regulation of plants.

368        To determine the biological function of DEGs during seed development, GO

369    classification analysis was carried out using Blast2GO. The results showed that 25 functional

370    groups, including 3 major ontologies, classified 63,192 GO terms annotated by the GO

371    database: biological process, cellular component, and molecular function. Many of these

16

372  DEGs were dominant catalytic activity, binding metabolism, cellular processes, cell parts and

373  cells (S5 Fig). In confirming whether there is specificity for development of seeds in relation

374  to their transcripts, orthologous *S. tora* genes were applied to gene ontology enrichment

375  analysis using the AgriGO program. In molecular function of GO ontologies, the level of

376  binding function was increased in the up-regulated DEGs. Among them, RNA binding

377  increased to a very high level. In addition, down-regulated DEGs showed an increase in the

378  catalytic activity function, and they also increased protein kinase activity, transferase activity,

379  and microtubule motor activity (S5 Fig).

380  To identify specific metabolic pathways that are responsible for the transcriptional

381  changes of enzymatic genes during seed development of *S. tora,* we performed MapMan

382  analysis with the expression data of genes showing at least 2-fold differential expression

383  between seed developmental stages. We made the figure to depict the biological processes of

384  interest, and display log2-normalized expression counts onto pictorial diagrams. Most of the

385  genes in cell metabolism are involved in cell wall metabolism, lipid metabolism,

386  carbohydrate metabolism, and secondary metabolism. The dynamic changes in metabolic

387  pathways during seed development were provided in Fig 5, in which we identified the

388  downward trend of overall transcription in the seed development process. In particular, it was

389  clear that lipid metabolism, precursor synthesis, flavonoid metabolisms, and

390  phenylpropanoids/phenolics metabolisms were down-regulated, while the FA synthesis of

391  lipid metabolism and the N-msc of secondary metabolism were up-regulated.

392

393  **Candidate gene families involved in anthraquinones biosynthesis**

394  *S. tora* is well known for its various therapeutic effects (e.g., for its anti-hypertensive,

395  diuretic, anti-cancer, anti-microbial and cholesterol-lowering effects). Each effect is caused

17

396    by various secondary metabolites produced in *S. tora*, the best known of these being

397    anthraquinone. The biosynthesis of anthraquinone shares isochorismate pathways with

398    phenylpropanoid and shares MEP/DOXP, MEV, and shikimate pathways with carotenoid and

399    flavonoid. In addition, the polyketide pathway is an important part of the anthraquinone

400    biosynthesis. To analyze the active biosynthesis of anthraquinones, we determined the

401    contents of seven compounds of the anthraquinone biosynthesis pathway in early- and late-

402    seed tissues. As seeds matured, anthraquinone compounds were more accumulated in late

403    seed than early seed (Fig 6 and Table 2). Among the seven compounds, gluco-obtusifolin has

404    the highest content in seed tissues (Fig 6 and Table 2). It is well known that aurantio-obtusin

405    is the most significant active compound [8] and is distributed mainly in the seed [44].

406    However, we found that low levels of aurantio-obtusin were observed at the early and late

407    developmental stages. A possible explanation for this reason is that aurantio-obtusin may

408    accumulate mainly in the matured and/or dry seed.

409        To observe gene expression levels of each parts and to compare the changes in gene

410    expression levels between different parts, their levels were normalized to the FPKM (reads

411    per kilobase of exon model per million mapped reads), and transcripts were hierarchically

412    clustered based on the Log2(FPKM+1), allowing us to observe the overall gene expression

413    pattern (Fig 7). In our study, there were 337 RNA-Seq and 212 Iso-Seq genes involved in *S.*

414    *tora* secondary metabolites, and they were classified into five pathways including the

415    MEP/DOXP, MEV, shikimate, carotenoid, and flavonoid/polyketide (Fig 7 and S5 Table).

416    There were 35 RNA-Seq and 24 Iso-Seq genes in *S. tora* for seven enzymes involved in

417    MEP/DOXP pathway and mevalonate pathway leading to production of precursor

418    dimethylallyl disphosphate (Fig 7 and S5 Table). They are also involved in the shikimate

419    pathway leading to the production of precursor 1,4-dihydroxy-2-napthoyl-CoA including 40

18

420    RNA-Seq and 31 Iso-Seq genes for 9 enzymes (DAHPS, DHQS, DHQD/SDH, SMK, EPSP,

421    CS, ICS, MenE, and MenB). In MEP/DOXP, 13 DXPS (1-deoxy-$_D$-xylulose-5-phosphate

422    synthase, EC 2.2.1.7) were expressed in anthraquinone synthesis. In them, DN49358_C0_g1

423    was expressed in large amounts up to the early stage of seed, but appeared to be greatly

424    reduced by the late stage. This gene was also expressed at high levels in leaf and root tissues.

425    Furthermore, DN27315_c0_g1 demonstrated higher levels of gene expression in leaf than in

426    other tissues. And only three of the 13 DXPS genes showed high levels of expression

427    independent of tissue and seed development. ISPD, CDPMEK, and ISPF genes were

428    identified in only 1 and 2, while HDS and HDR were identified in more frequent. HDS and

429    HDR were identified in genes 8 and 6, and HDS ((E)-4-hydroxy-3-methylbut-2-enyl-

430    diphosphate synthase, DN48094_c1_g1) and HDR (4-hydroxy-3-methylbut-2-enyl-

431    diphosphate reductase, DN25595_c0_g1) showed high levels of expression regardless of

432    tissue and seed development. In the MEV pathway, ACCA (acetyl-CoA carboxylase) was

433    identified in 29 genes, and 3 genes (DN51063_c1_g1, DN51063_c2_g1, and

434    DN72707_c0_g1) sustained high levels of expression independent of tissue and seed

435    development. Conversely, one HMGR (DN9882_c0_g1) was down-stream of expression

436    level. Except for some genes, ACCA, HMGS, HMGR, MK, PMK, and MPD of expression

437    levels are down-stream, and 1 of 4 IPPS (isopentenyl-diphosphate delta-isomerase,

438    DN67602_c1_g1) genes showed high level of expression independent of tissue and seed

439    development.

440        Anthraquinones are also known to be produced from acetyl-CoA and malonyl-CoA

441    through polyketide pathway in plants. Chalcone synthase (CHS), a type III polyketide

442    synthase, is an important enzyme involved in the polyketide pathway [45]. We have

443    identified 27 RNA-Seq and 23 Iso-Seq genes encoding for enzyme involved in type III

19

444    polyketide synthase (S5 Table). As a ubiquitous enzyme in higher plants, CHS is known to

445    produce flavonoids by catalyzing the sequential decarboxylative reaction with 3 malonyl-

446    CoA and p-coumaroyl-CoA as a starter and extender unit, respectively [46]. It was also

447    suggested that polyketide synthase could form an anthraquinone precursor using acetyl-CoA

448    and malonyl-CoA. And the formed precursor, octaketide is cyclized by PKC-encoding

449    polyketide cyclase, and usually forms three-ring structures named A, B, and C rings [47]. The

450    formed intermediate is modified by P450 to produce anthraquinone or emodin anthrone, and

451    also to produce sennoside by modification of glycosyltransferases. These 27 PKS gene sizes

452    averaged 584.03 bp, and the longest was 1,580 bp. Among them, only 3 genes

453    (DN50459_c0_g1, DN2403_c0_g2, and DN50459_c0_g2) showed high levels of expression

454    change in seed development. It seems that these genes are changing a lot in order to make the

455    backbones of the flavonoid and carotenoid components needed for survival in the later stages

456    of seed development. In particular, 5 genes (DN17347_c0_g1, DN50624_c4_g3,

457    DN69520_c0_g1, DN50624_c4_g1, and DN50624_c4_g4) showed a large amount of

458    expression in the early part of the seed, whereas in the latter part, the level of expression

459    decreased sharply, suggesting that those genes play a very important role in the biosynthesis

460    of the backbone of the material needed in early seed development.

461         In general, glycosylation is carried out at the end of secondary metabolites

462    biosynthesis and improve the solubility and stability of the secondary metabolites. In nature,

463    UDP-glycosyltransferases (UGT) normally facilitates glycosylation, and makes the natural

464    product with glucose at the hydroxyl group [48]. In our study, there were 59 genes in seed

465    stage of *S. tora*. Based on the results, 33 out of 59 genes showed more expression at the late-

466    seed than at the early-seed stage, whereas 26 showed more expression at the early-seed stage

467    (Fig 7 and S5 Table). The degree of expression of the seven genes (DN131354_c0_g1,

20

468     DN67413_c0_g1, DN49988_c0_g2, DN50503_c0_g2, DN82643_c0_g1, DN17331_c0_g2,

469     and DN137099_c0_g1) seems to increase rapidly during the growth of the seed, which seems

470     to be necessary for the process of stockpiling the energy required for seed germination. In

471     addition, DN17331_c0_g2 and DN82643_c0_g1 seem to have a great effect on the

472     glycosylation during seed development because they undergo a significant amount of change.

473     Conversely, the expression level of the four genes (DN50189_c2_g1, DN11235_c0_g1,

474     DN62590_c0_g1, and DN76515_c0_g1) seemed to decrease rapidly, and the remaining 22

475     genes were found to be expressed with a relatively small decrease.

476

# Conflict of interest

478     The authors declare that they have no conflict of interest.

479

# Acknowledgements

485

# Author Contributions

487     **Funding acquisition:** Sang-Ho Kang

488     **Data curation:** Joon-Soo Sim, Chang-Muk Lee, So-Ra Han

489     **Methodology:** So Youn Won, Soo-Jin Kwon, Jung Sun Kim

490     **Writing – original draft:** Sang-Ho Kang, Woo-Haeng Lee

491    **Writing – review & editing:** Sang-Ho Kang, Chang-Kug Kim, Tae-Jin Oh

492

# References

494    1.    Jain R, Sharma P, Jain SC. Chemical analysis of the roots of *Cassia tora*. Asian J

495          Chem. 2010;22(10):7585-90.

496    2.    Patil UK, Saraf S, Dixit VK. Hypolipidemic activity of seeds of *Cassia tora* Linn. J

497          Ethnopharmacol. 2004;90(2-3):249-52.

498    3.    Pawar HA, D'mello PM. *Cassia tora* Linn.: An overview. Int J Pharmaceut Sci Res.

499          2011;2(9):2286-91.

500    4.    Zhao X, Wang Q, Qian Y, Pang L. *Cassia tora* L. (Jue-ming-zi) has anticancer

501          activity in TCA8113 cells *in vitro* and exerts anti-metastatic effects *in vivo*. Oncol

502          Lett. 2013;5(3):1036-42.

503    5.    Jang DS, Lee GY, Kim YS, Lee YM, Kim C-S, Yoo JL, et al. Anthraquinones from

504          the seeds of *Cassia tora* with inhibitory activity on protein glycation and aldose

505          reductase. Biol Pharm Bull. 2007;30(11):2207-10.

506    6.    Shukla SK, Kumar A, Terrence M, Yusuf J, Singh VP, Mishra M. The probable

507          medicinal usage of *Cassia tora*: An overview. OnLine J Biol Sci. 2013;13(1):13-7.

508    7.    Jain S, Patil UK. Phytochemical and pharmacological profile of Cassia tora Linn. -

509          An overview. Indian J Nat Prod Resour. 2010;1(4):430-7.

510    8.    Kim M, Lim SJ, Lee HJ, Nho CW. *Cassia tora* seed extract and its active compound

511          aurantio-obtusin inhibit allergic responses in IgE-mediated mast cells and

512          anaphylactic models. J Agric Food Chem. 2015;63(41):9037-46.

513    9.    Chien S-C, Wu Y-C, Chen Z-W, Yang W-C. Naturally occurring anthraquinones:

514          Chemistry and therapeutic potential in autoimmune diabetes. Evidence-Based

515        Complementary and Alternative Medicine. 2015;2015:1-13.

516  10.    Kim Y-M, Lee C-H, Kim H-G, Lee H-S. Anthraquinones isolated from *Cassia tora*

517        (Leguminosae) seed show an antifungal property against phytopathogenic fungi. J

518        Agric Food Chem. 2004;52:6096-100.

519  11.    Azelmat J, Larente JF, Grenier D. The anthraquinone rhein exhibits synergistic

520        antibacterial activity in association with metronidazole or natural compounds and

521        attenuates virulence gene expression in *Porphyromonas gingivalis*. Arch Oral Biol.

522        2015;60(2):342-6. Epub 2014/12/03.

523  12.    Panigrahi GK, Ratnasekhar CH, Mudiam MKR, Vashishtha VM, Raisuddin S, Das

524        M. Activity-guided chemo toxic profiling of *Cassia occidentalis* (CO) seeds:

525        Detection of toxic compounds in body fluids of CO-exposed patients and

526        experimental rats. Chem Res Toxicol. 2015;28:1120-32.

527  13.    Panigrahi GK, Yadav A, Mandal P, Tripathi A, Das M. Immunomodulatory potential

528        of rhein, an anthraquinone moiety of *Cassia occidentalis* seeds. Toxicol Lett.

529        2016;245:15-23.

530  14.    Wong S-M, Wong MM, Seligmann O, Wagner H. New antihepatotoxic naphtho-

531        pyrone glycosides from the seeds of *Cassia tora*. Plant Med. 1989;55:276-80.

532  15.    Yen G-C, Chung D-Y. Antioxidant effects of extracts from *Cassia tora* L. prepared

533        under different degrees of roasting on the oxidative damage to biomolecules. J Agric

534        Food Chem. 1999;47:1326-32.

535  16.    Hyun SK, Lee H, Kang SS, Chung HY, Choi JS. Inhibitory activities of *Cassia tora*

536        and its anthraquinone constituents on angiotensin-converting enzyme. Phytother Res.

537        2009;23(2):178-84.

538  17.    Baba SA, Mohiuddin T, Basu S, Swarnkar MK, Malik AH, Wani ZA, et al.

539    Comprehensive transcriptome analysis of *Crocus sativus* for discovery and

540    expression of genes involved in apocarotenoid biosynthesis. BMC Genomics.

541    2015;16:698.

542  18.  D'Agostino N, Pizzichini D, Chiusano ML, Giuliano G. An EST database from

543    saffron stigmas. BMC Plant Biol. 2007;7:53.

544  19.  Jain M. Next-generation sequencing technologies for gene expression profiling in

545    plants. Brief Funct Genomics. 2012;11(1):63-70.

546  20.  Jo IH, Lee J, Hong CE, Lee DJ, Bae W, Park SG, et al. Isoform sequencing provides

547    a more comprehensive view of the *Panax ginseng* transcriptome. Genes (Basel).

548    2017;8(9).

549  21.  Li J, Harata-Lee Y, Denton MD, Feng Q, Rathjen JR, Qu Z, et al. Long read

550    reference genome-free reconstruction of a full-length transcriptome from *Astragalus*

551    *membranaceus* reveals transcript variants involved in bioactive compound

552    biosynthesis. Cell Discov. 2017;3:17031.

553  22.  Gonzalez-Garay ML. Transcriptomics and gene regulation: Introduction to isoform

554    sequencing using pacific biosciences technology (Iso-Seq): Springer; 2016.

555  23.  Pouladi N, Achour I, Li H, Berghout J, Kenost C, Gonzalez-Garay ML, et al.

556    Biomechanisms of comorbidity: Reviewing integrative analyses of multi-omics

557    datasets and electronic health records. Yearb Med Inform. 2016;(1):194-206.

558  24.  Patel RK, Jain M. NGS QC Toolkit: a toolkit for quality control of next generation

559    sequencing data. PLoS One. 2012;7(2):e30619.

560  25.  Schmieder R, Lim YW, Edwards R. Identification and removal of ribosomal RNA

561    sequences from metatranscriptomes. Bioinformatics. 2012;28(3):433-5.

562  26.  Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. *De*

24

563     *novo* transcript sequence reconstruction from RNA-seq using the Trinity platform for

564     reference generation and analysis. Nat Protoc. 2013;8(8):1494-512.

565     27.    Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of

566     protein or nucleotide sequences. Bioinformatics. 2006;22(13):1658-9.

567     28.    Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat Methods.

568     2012;9(4):357-9.

569     29.    Li B, Dewey CN. RSEM: accurate transcript quantification from RNA-Seq data with

570     or without a reference genome. BMC Bioinformatics. 2011;12:323.

571     30.    Anders S, Huber W. Differential expression analysis for sequence count data.

572     Genome Biol. 2010;11:R106.

573     31.    Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment

574     search tool. J Mol Biol. 1990;215:403-10.

575     32.    Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene

576     ontology: tool for the unification of biology. The Gene Ontology Consortium. Nat

577     Genet. 2000;25(1):25-9.

578     33.    Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. Blast2GO: a

579     universal tool for annotation, visualization and analysis in functional genomics

580     research. Bioinformatics. 2005;21(18):3674-6.

581     34.    Du Z, Zhou X, Ling Y, Zhang Z, Su Z. agriGO: a GO analysis toolkit for the

582     agricultural community. Nucleic Acids Res. 2010;38(Web Server issue):W64-70.

583     35.    Livak KJ, Schmittgen TD. Analysis of relative gene expression data using real-time

584     quantitative PCR and the 2(-Delta Delta C(T)) Method. Methods. 2001;25(4):402-8.

585     36.    Yang CQ, Fang X, Wu XM, Mao YB, Wang LJ, Chen XY. Transcriptional

586     regulation of plant secondary metabolism. J Integr Plant Biol. 2012;54(10):703-12.

587  37.  Omidbakhshfard MA, Proost S, Fujikura U, Mueller-Roeber B. Growth-regulating

588       factors (GRFs): A small transcription factor family with important functions in plant

589       biology. Mol Plant. 2015;8(7):998-1010.

590  38.  Liu J, Hua W, Yang H-L, Zhan G-M, Li R-J, Deng L-B, et al. The BnGRF2 gene

591       (GRF2-like gene from Brassica napus) enhances seed oil production through

592       regulating cell number and plant photosynthesis. J Exp Bot. 2012;63(10):3727-40.

593  39.  Zhang D-F, Li B, Jia G-Q, Zhang T-F, Dai J-R, Li J-S, et al. Isolation and

594       characterization of genes encoding GRF transcription factors and GIF transcriptional

595       coactivators in maize (*Zea mays* L.). Plant Science. 2008;175(6):809-17.

596  40.  Link M, Rausch T, Greiner S. In *Arabidopsis thaliana*, the invertase inhibitors

597       AtC/VIF1 and 2 exhibit distinct target enzyme specificities and expression profiles.

598       FEBS Lett. 2004;573(1-3):105-9.

599  41.  Kang SH, Lee JY, Lee TH, Park SY, Kim CK. *De novo* transcriptome assembly of

600       the Chinese pearl barley, adlay, by full-length isoform and short-read RNA

601       sequencing. PLoS One. 2018;13(12):e0208344.

602  42.  Hansen CH, Du L, Naur P, Olsen CE, Axelsen KB, Hick AJ, et al. CYP83b1 is the

603       oxime-metabolizing enzyme in the glucosinolate pathway in Arabidopsis. J Biol

604       Chem. 2001;276(27):24790-6.

605  43.  Kufner I, Koch W. Stress regulated members of the plant organic cation transporter

606       family are localized to the vacuolar membrane. BMC Res Notes. 2008;1:43.

607  44.  Deng Y, Zheng H, Yan Z, Liao D, Li C, Zhou J, et al. Full-length transcriptome

608       survey and expression analysis of *Cassia obtusifolia* to discover putative genes

609       related to aurantio-obtusin biosynthesis, seed formation and development, and stress

610       response. Int J Mol Sci. 2018;19(9).

611  45.  Pandith SA, Dhar N, Rana S, Bhat WW, Kushwaha M, Gupta AP, et al. Functional

612        promiscuity of two divergent paralogs of type III plant polyketide synthases. Plant

613        Physiol. 2016;171(4):2599-619.

614  46.  Austin MB, Noel JP. The chalcone synthase superfamily of type III polyketide

615        synthases. Nat Prod Rep. 2003;20:79-110.

616  47.  Rama Reddy NR, Mehta RH, Soni PH, Makasana J, Gajbhiye NA, Ponnuchamy M,

617        et al. Next generation sequencing and transcriptome analysis predicts biosynthetic

618        pathway of sennosides from *Senna* (*Cassia angustifolia* Vahl.), a non-model plant

619        with potent laxative properties. PLoS One. 2015;10(6):e0129422.

620  48.  Ross J, Li Y, Lim E-K, Bowles DJ. Higher plant glycosyltransferases. Genome Biol.

621        2001;2(2):3004.1-.6.

622

27

623     **Table 1.** Assembly statistics of the *S. tora* transcriptome by RNA-Seq and Iso-Seq.

| Assembly statistics | RNA-Seq | Iso-Seq |
|---|---|---|
| Number of unigenes | 118,635 | 39,364 |
| Total size (bp) | 98,734,027 | 112,216,332 |
| Minimum length (bp) | 300 | 435 |
| Maximum length (bp) | 18,622 | 6,814 |
| Average length (bp) | 832 | 2,851 |
| N50 length (bp) | 1,082 | 3,513 |
| GC contents (%) | 39.51 | 38.60 |

**Table 2.** Anthraquinone contents in the early and late seeds.

| Compounds | Formula | RT[a] | Contents (ug/g) | |
|---|---|---|---|---|
| | | | **Early Seed** | **Late Seed** |
| Gluco-obtusifolin | $C_{22}H_{22}O_{10}$ | 5.31 | 80.72[b] | 141.27 |
| Aurantio-obtusin | $C_{17}H_{14}O_7$ | 5.27 | 1.07 | 1.08 |
| Chryso-obtusin | $C_{19}H_{18}O_7$ | 7.68 | 1.39 | 0.69 |
| Obtusin | $C_{18}H_{16}O_7$ | 8.07 | 0.83 | 0.60 |
| Obtusifolin | $C_{16}H_{12}O_5$ | 8.19 | 1.38 | 0.26 |
| Chrysophanol | $C_{15}H_{10}O_4$ | 10.65 | 11.39 | 6.69 |
| Physcion | $C_{16}H_{12}O_5$ | 11.14 | 2.02 | 0.96 |
| Total | - | - | 98.80 | 151.55 |

[a] indicates retention time.

[b] represents mean of three replicate experiments.

29

## Figure legends

**Fig 1. The length distribution of transcripts in *S. tora*.** X and Y axis represent unigene lengths and percent of unigene length distribution, respectively.

**Fig 2. Histogram of gene ontology (GO) classification from RNA-Seq and Iso-Seq.** The results are summarized in three main categories: biological process, molecular function, and cellular component.

**Fig 3. Heat map of top 50 differentially expressed genes between early- and late-stages of seed development in *S. tora*.** Heatmap showing differentially expressed genes between early and late stages of seed development in *S. tora*. Color scale representing normalized expression values is shown at the bottom.

**Fig 4. Real-time PCR validation of gene expression obtained via RNA-Seq.** All the real-time PCR experiments were performed at least three times in each independent biological experiment (3 replicates). Error bars represent SEM from triplicates.

**Fig 5. MapMan metabolism overview maps showing differences in transcript levels during seed development.** MapMan software was used to provide a snapshot of modulated genes over the main metabolic pathways. Log2 fold changes values are represented. Up-regulated and down-regulated transcripts are shown in red and blue, respectively.

**Fig 6. GC-MS analysis of anthraquinone during seed development.** Seven anthraquinone

653   levels in the early seed (A) and in the late seed (B).

654

655   **Fig 7. The up-down of putative genes of anthraquinone-biosynthetic pathway in *S. tora.***

656   It was normalized to the FPKM to compare the changes in gene expression levels between

657   different parts of *S. tora*. Total gene expression levels were clustered based on the Log2

658   (FPKM +1). DXPS, 1-Deoxy-$_D$-xylulose-5-phosphate synthase (EC 2.2.1.7); DXR, 1-Deoxy-

659   $_D$-xylulose-5-phosphate reductoisomerase (EC 1.1.1.267); ISPD, 2-C-Methyl-$_D$-erythritol 4-

660   phosphate cytidylyltransferase (EC 2.7.7.60); CDPMEK, 4-Diphosphocytidyl-2-C-methyl-$_D$-

661   erythritol kinase (EC 2.7.1.148); ISPF, 2-C-Methyl-$_D$-erythritol 2,4-cyclodiphosphate

662   Synthase (EC 4.6.1.12); HDS, (E)-4-Hydroxy-3-methylbut-2-enyl-diphosphate synthase (EC

663   1.17.7.1); HDR, 4-Hydroxy-3-methylbut-2-enyl diphosphate reductase (EC 1.17.1.2); ACCA,

664   Acetyl-CoA carboxylase (EC 6.4.1.2); HMGS, Hydroxymethylglutaryl-CoA synthase (EC

665   2.3.3.10); HMGR, Hydroxymethylglutaryl-CoA reductase (EC 1.1.1.34); MK, Mevalonate

666   kinase (EC 2.7.1.36); PMK, Phosphomevalonate kinase (EC 2.7.4.2); MPD, Methyl parathion

667   hydrolase (EC 3.1.8.1); IPPS, Isopentenyl-diphosphate delta-isomerase (EC 5.3.3.2);

668   DAHPS, 3-Deoxy-7-phosphoheptulonate synthase (EC 2.5.1.54); DHQS, 3-Dehydroquinate

669   synthase (EC 4.2.3.4); DHQD/SDH, 3-Dehydroquinate dehydratase/shikimate dehydrogenase

670   (EC 4.2.1.10/1.1.1.25); SMK, Shikimate kinase (EC 2.7.1.71); EPSP, 3-Phosphoshikimate 1-

671   carboxyvinyltransferase (EC 2.5.1.19); CS, Chorismate synthase (EC 4.2.3.5); ICS,

672   Isochorismate synthase (EC 5.4.4.2); PHYLLO, 2-Succinyl-5-enolpyruvyl-6-hydroxy-3-

673   cyclohexene-1-carboxylic acid synthase (EC 2.2.1.9); MenE, 2-Succinylbenzoate-CoA ligase

674   (EC 6.2.1.26); MenB, 1,4-Dihydroxy-2-naphthoyl-CoA synthase (EC 4.1.3.36); GGPS,

675   Geranylgeranyl diphosphate synthase (EC 2.5.1.1); PSY, Phytoene synthase (EC 2.5.1.32);

676   PDS, Phytoene desaturase (EC 1.3.99.30); ZDS, Zeta-carotene desaturase (EC 1.3.5.6);

31

677  LYCB, Lycopene beta-cyclase (EC 5.5.1.19); LYCE, Lycopene epsilon-cyclase (EC

678  5.5.1.18); BCH, Beta-carotene hydroxylase (EC 1.14.13.129); ZEP, Zeaxanthin epoxidase

679  (EC 1.14.15.21); PAL, Phenylalanine ammonia-lyase (EC 4.3.1.24); C4H, Cinnamate-4-

680  hydroxylase (EC 1.14.13.11); 4CL, 4-Coumarate-CoA ligase (EC 6.2.1.12); and CHS,

681  Chalcone synthase (EC 2.3.1.74).

682

# Supporting information

**S1 Table. Gene-specific primers used for tissue-specific qRT-PCR.**

**S2 Table. General properties of the reads produced by Illumina Hiseq 2500 sequencing platform.**

**S3 Table. General properties of the reads produced by PacBio sequencing platform.**

**S4 Table. Tissue-enriched and specific transcription factors (TFs) distribution of each tissue.**

**S5 Table. Gene associated with the secondary metabolite pathway in *S. tora.***

**S1 Fig. The distribution of annotated unigenes by various public protein databases.** Venn diagram showing the proportion of annotated unigenes in NCBI Nr, KEGG, Swiss-Prot, and GO databases with RNA-Seq (**A**) and Iso-Seq (**B**).

**S2 Fig. Species distribution of the top BLAST hits.** Top-hit species from RNA-Seq and Iso-Seq were calculated based on sequence alignments with the lowest E-value obtained from BLAST.

**S3 Fig. Distribution of TF families of *S. tora.*** Distribution of transcripts (3,284 for RNA-Seq and 3,576 for Iso-Seq) that encode for transcription factors (**A**). Number of transcripts exhibiting specific expression in different tissues has been indicated by bar and table (**B**). Tissue-specific shows 10-fold higher FPKM in one tissue compared with three tissues, and tissue-enriched represents 5-fold higher FPKM compared with other tissues.

**S4 Fig. Heatmaps representing the top 10 genes that showed tissue-specific expression in the *S. tora* leaf, root, and early and late seeds.** Red represents high abundance and green represents low abundance.

**S5 Fig. AgriGo analysis of upregulated and downregulated genes during seed**

33

706    **development.** A total of 4,935 (up-regulated, **A**) and 9,890 (down-regulated, **B**) genes with

707    Molecular terms are represented by increasingly red colors. GO term enrichment was

708    performed    using    single    enrichment    analysis    (SEA)    tool    on    AgriGo

709    (http://bioinfo.cau.edu.cn/agrigo/). Box colors indicates levels of statistical significance:

710    yellow=0.05; orange=e-05; and red=e-09.

Figure 1

Figure 2

Oleosin Bn-III
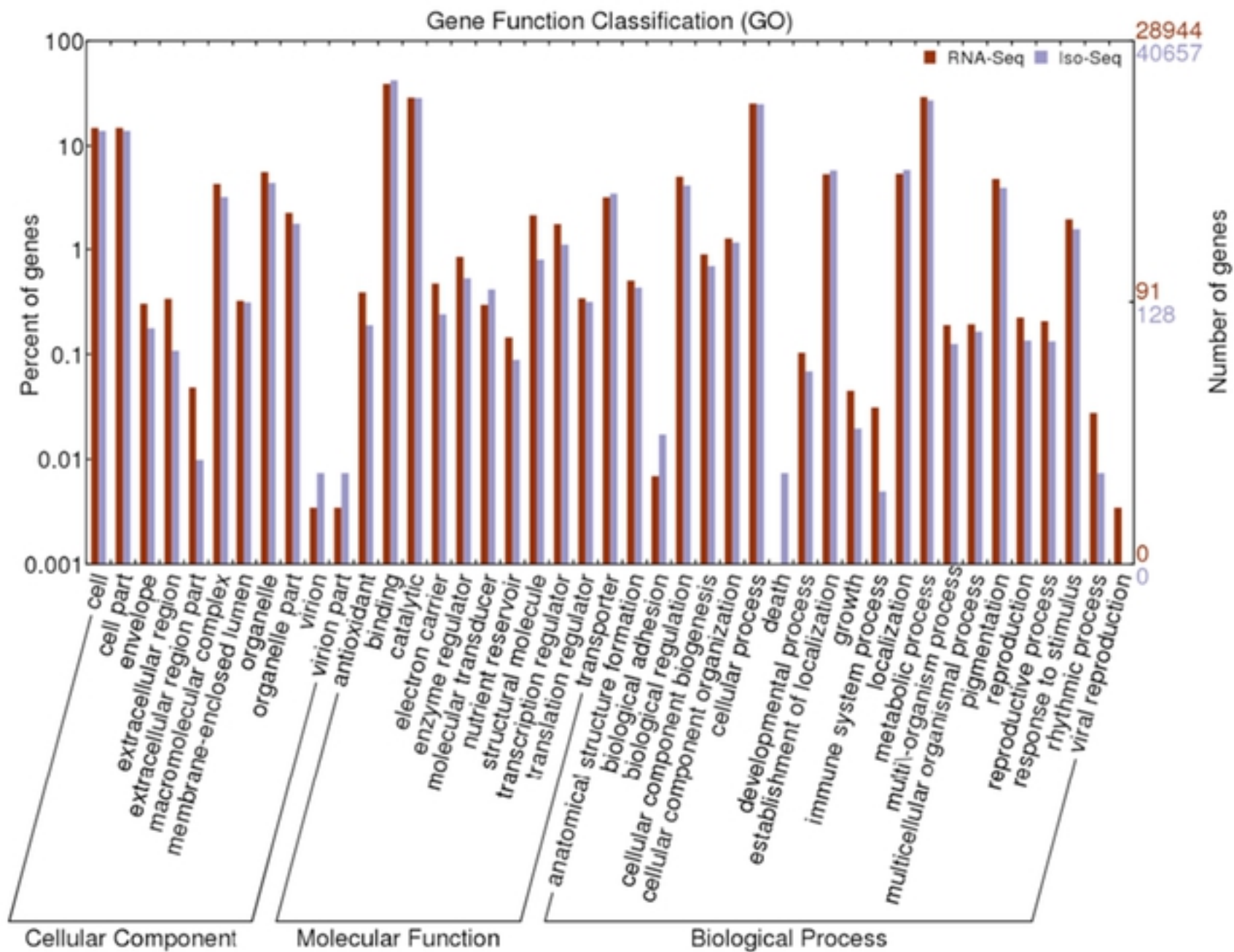Vicilin-like antimicrobial peptides 2-1
Globulin-1 S allele
Metallothionein-like protein type 2
1-aminocyclopropane-1-carboxylate oxidase
Oleosin 21.2 kDa
Unknown
Unknown
Snakin-2
Beta-glucosidase 44
Unknown
Chalcone synthase 2
Chalcone synthase
Unknown
Isoliquiritigenin 2'-O-methyltransferase
RING-H2 finger protein ATL72-like
Peroxidase 42
Cell wall / vacuolar inhibitor of fructosidase 2
Unknown
MLP-like protein 43
18 kDa seed maturation protein
Aldose reductase
Embryonic protein DC-8
Embryonic protein DC-8
Desiccation protectant protein Lea14 homolog
Defensin-like protein
Reticulon-like protein B13
Low-temperature-induced 65 kDa protein
Seed maturation PM36-like protein
Late embryogenesis abundant protein
Seed biotin-containing protein SBP65
Late embryogenesis abundant protein
Seed biotin-containing protein SBP65
Protein SLE1
Unknown
Unknown
17.3 kDa class I heat shock protein
Desiccation-related protein PCC13-62
Late embryogenesis abundant protein
Late embryogenesis abundant protein
Protein SLE2
dehydrin Rab18-like
Unknown
Oleosin 18.5 kDa
Glycinin G3
P24 oleosin isoform B
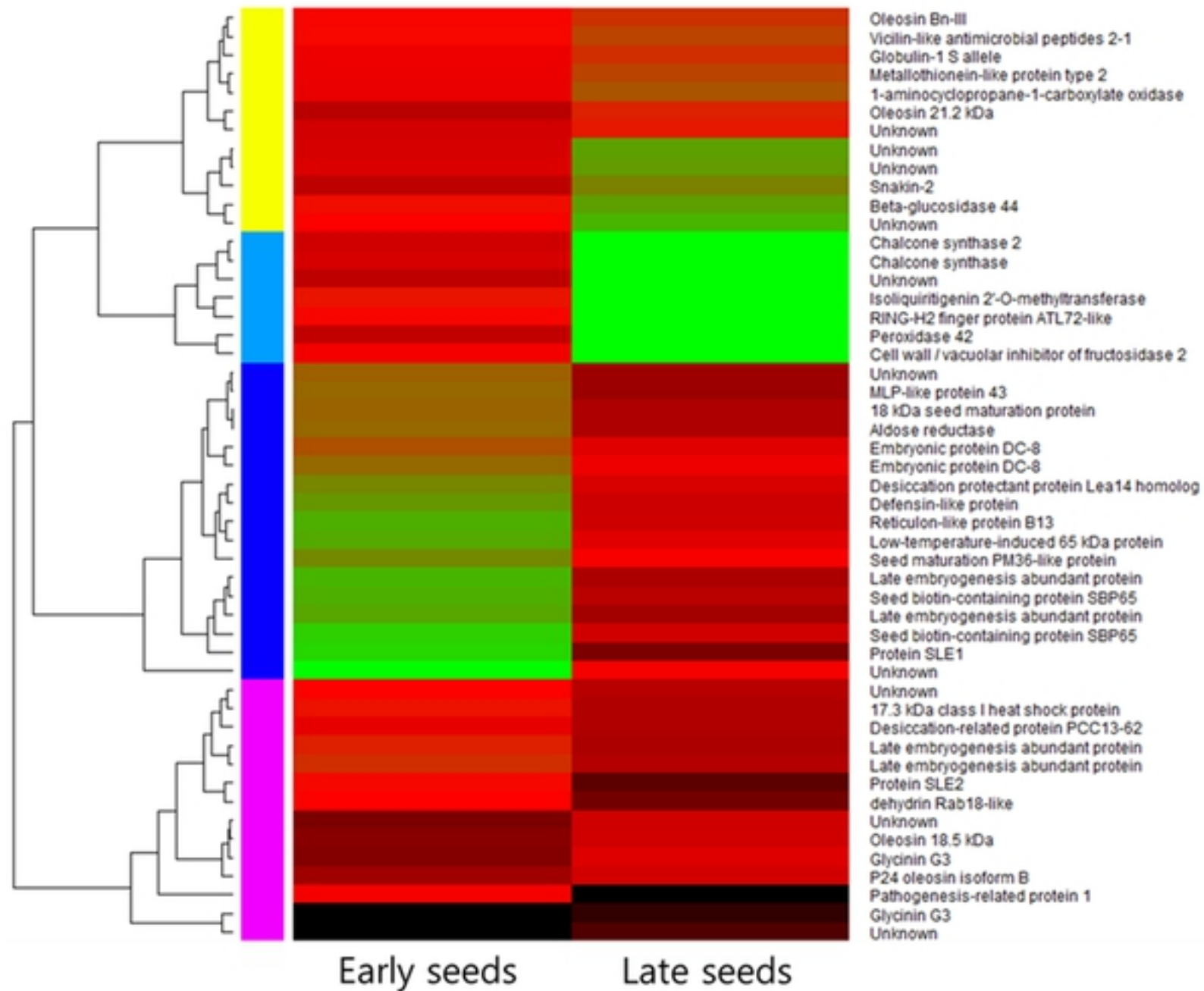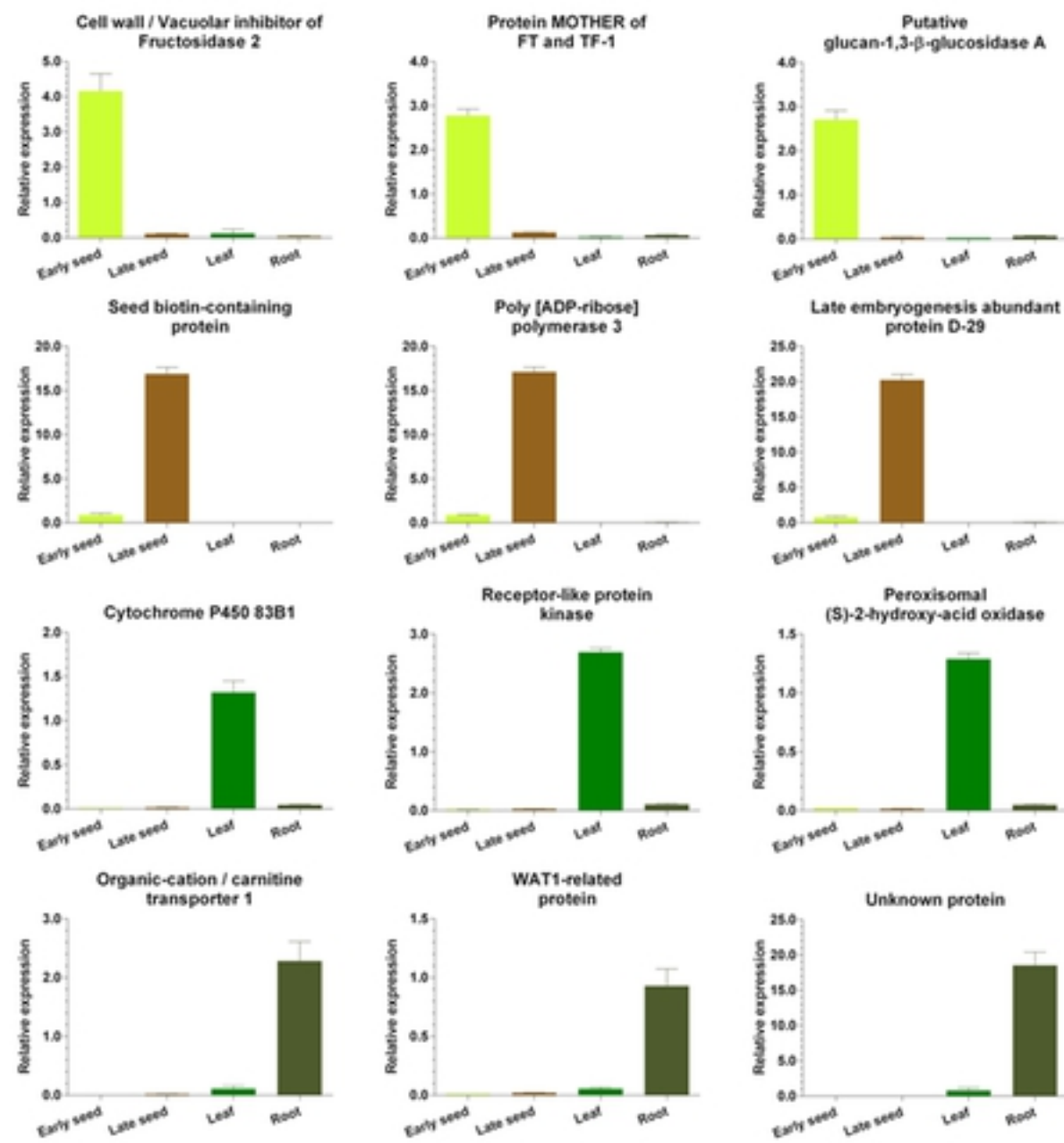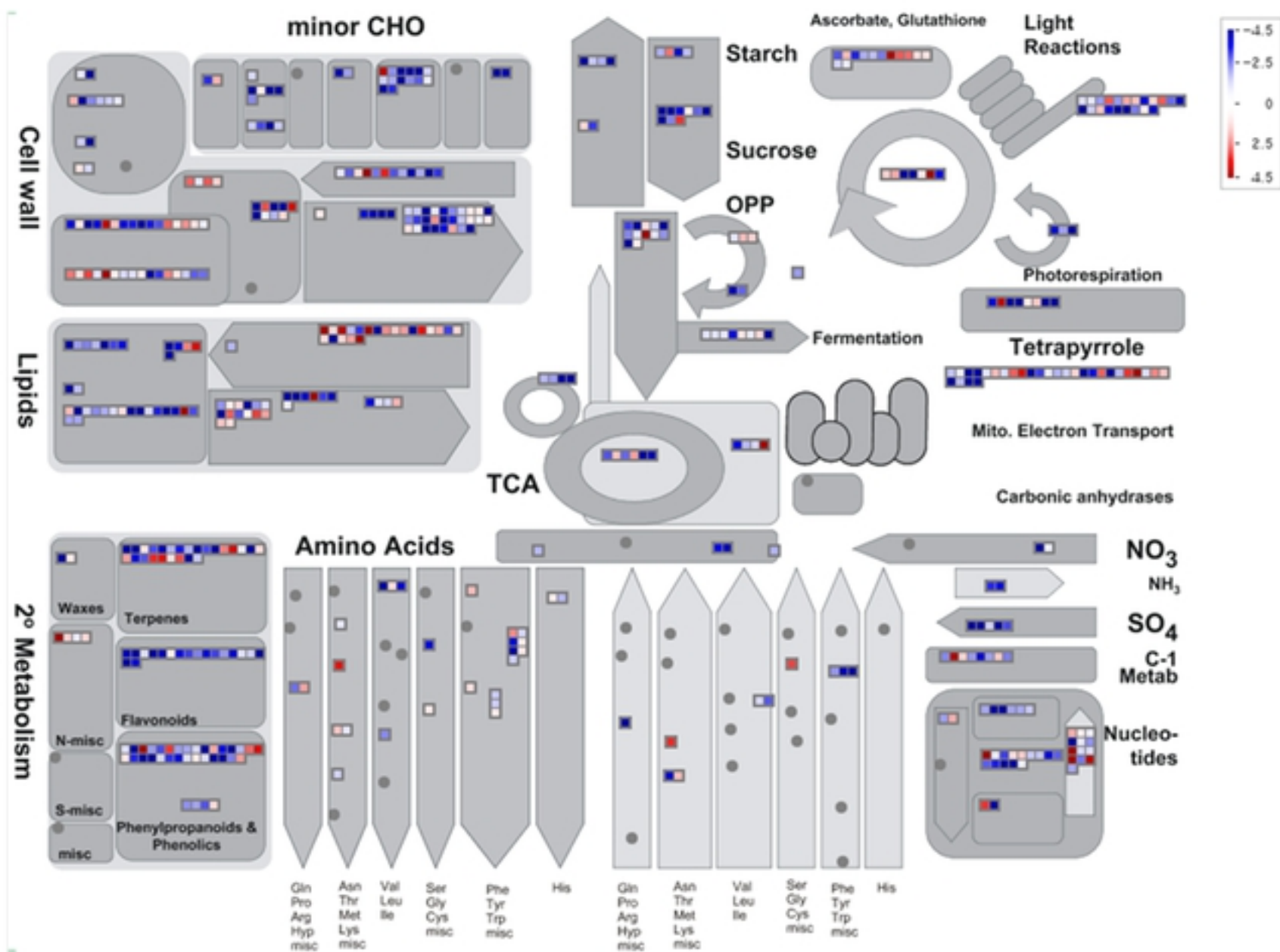Pathogenesis-related protein 1
Glycinin G3
Unknown

Early seeds          Late seeds

Figure 3

Figure 4

Figure 5

Figure 6

Figure7