

1 **Title: Comparing Time Series Transcriptome Data Between Plants Using A Network Module**

2 **Finding Algorithm**

3

4 Jiyoung Lee^{1,3}, Lenwood S. Heath², Ruth Grene³, Song Li^{1,3}

5

6 1. Ph.D. program in Genetics, Bioinformatics and Computational Biology, Virginia Polytechnic Institute
7 and State University, Blacksburg, VA, 24061.

8

9 2. Department of Computer Science, Virginia Polytechnic Institute and State University, Blacksburg, VA,
10 24061.

11

12 3. School of Plant and Environmental Sciences, Virginia Polytechnic Institute and State University,
13 Blacksburg, VA, 24061.

14

15

16

17

18 **Author emails:**

19 Jiyoung Lee: jylee43@vt.edu

20 Ruth Grene: grene@vt.edu

21 Lenwood S. Heath: heath@vt.edu

22 Song Li: songli@vt.edu

23

24 **Corresponding authors:** Song Li, songli@vt.edu, Phone: 540-231-2756

25

26

27 **Running Head: Comparative Transcriptome Analysis**

28 **ABSTRACT**

29 Comparative transcriptome analysis is the comparison of expression patterns between homologous genes
30 in different species. Since most molecular mechanistic studies in plants have been performed in model
31 species including Arabidopsis and rice, comparative transcriptome analysis is particularly important for
32 functional annotation of genes in other plant species. Many biological processes, such as embryo
33 development, are highly conserved between different plant species. The challenge is to establish one-to-
34 one mapping of the developmental stages between two species. In this protocol, we solve this problem by
35 converting the gene expression patterns into a co-expression network and then apply network module-
36 finding algorithms to the cross-species co-expression network. We describe how to perform such analysis
37 using bash scripts for preliminary data processing and R programming language, which implemented
38 simulated annealing method for module finding. We also provide instructions on how to visualize the
39 resulting co-expression networks across species.

40

41 **Keywords**

42 Comparative transcriptome analysis, Network, Sequence homology, Arabidopsis, Soybean, Emybro
43 development

44

45 **INTRODUCTION**

46 Expression analysis is commonly used to understand the tissue or stress specificity of genes in
47 large gene families [1–5]. The goal of comparative transcriptome analysis is to identify conserved co-
48 expressed genes in two or more species [3,6,7]. The traditional definition of orthologous genes is based
49 solely on sequence homology [8–11] and syntenic relationships [2,12–14] and not on gene expression
50 patterns. In contrast, comparative transcriptome analysis combines a comparison of gene sequences with a
51 comparison of expression patterns between homologous genes in different species. Homologous genes
52 have been reported to be expressed either at different developmental stages, in different tissue types,

53 and/or under different stress conditions [3,15–17]. This documented divergence of expression patterns
54 provides crucial evidence for the existence of functional divergence of homologous genes across species
55 [18,19]. Therefore, comparative transcriptome analysis is an important tool for distinguishing those genes
56 that have retained functional conservation from those that have undergone functional divergence.
57 Comparative transcriptome analysis is particularly important for plant research, since most molecular
58 mechanistic studies in plants have been performed in model species, primarily Arabidopsis [20]. The
59 consequence of this narrow focus is that the functional annotation of the genes of many other plant
60 species relies solely on sequence comparisons with Arabidopsis [21].

61
62 To compare transcriptomes between any two species, a first step is to establish homologous
63 relationships between proteins in the two species. A second step is to identify expression data obtained
64 from experiments that are performed under similar conditions or tissue types. The third step is to compare
65 the expression patterns between the two data sets. In this protocol, we will compare published time course
66 seed embryo expression data from Arabidopsis [22] with data from the same tissue in soybean [23] as a
67 demonstration of how to apply computational tools to comparative transcriptome analysis.

68
69 In contrast with the time course data examined here, many other datasets have been reported from
70 “treatment-control” experiments (one time point only, two treatment conditions). For example, soybean
71 roots were treated with drought stress in one experiment [4]. To address the question of functional
72 conservation versus functional divergence within gene families, these soybean root data can be compared
73 with transcriptome data from Arabidopsis roots, under a similar stress [24]. This is a relatively simple
74 problem, because, in both experiments, we can identify lists of differentially expressed genes in response
75 to the same or similar treatments. It is a simple two-step process to identify conserved co-expressed genes
76 for treatment-control experiments. First, one needs to identify a list of gene pairs that are homologous
77 between these two species. A simple BLAST search or other more sophisticated approaches such as OMA,
78 EggNog, or Plaza [9,10,12] can be used to identify homologous genes. Second, the two lists of

79 differentially expressed genes can be compared to find whether any pairs of these homologous genes
80 appear in both lists.

81
82 In this article, we are focusing on a more complex scenario: two time-series experiments were
83 performed for the same developmental process in two different species [25]. Time course data provide
84 more data points than simple treatment-control experiments and, thus, can reveal relationships based on
85 development between homologous genes in two organisms. However, this is also challenging, because
86 the number of time points in the two experiments are different. It can be challenging to precisely match
87 developmental stages between two species, although some excellent approaches have been proposed
88 [25,26]. Despite the difficulty of establishing one-to-one mapping between the developmental stages of
89 two species, many biological processes, such as embryo development, are known to be highly conserved
90 between different plant species that are compared in comparative transcriptome analysis [27,28]. One way
91 to solve this developmental stage problem is to convert the gene expression patterns into a co-expression
92 network and then apply network alignment or network module-finding algorithms to these co-expression
93 networks [29]. Transforming expression data to a network form simplifies the problem and allows
94 exploration using well established network algorithms [30,31]. In this protocol, we describe how to
95 perform such analysis using a published simulated annealing method [29]. We also discuss how to
96 visualize the resulting co-expression networks across species [32] and the results from different choices of
97 homology finding methods.

98

99 **2. Install software and download experimental data**

100 All scripts used in this analysis can be obtained from github using the following command (Note 4.1).

101 The “\$” means the command is executed under a Linux terminal (Note 4.2).

102

103 \$ git clone <https://github.com/LiLabAtVT/CompareTranscriptome.git> ATH_GMA

104

105 You can replace “ATH_GMA” with another folder name that better represents your project. All scripts in
106 this project are tested under the project folder created by the “git clone” command (default ATH_GMA).

107

108 *Necessary Resources*

109 This protocol was tested under CentOS 7, which is a Linux operating system. The steps described in this
110 protocol can be used in most UNIX compatible operating systems; this includes all major Linux
111 distributions, and Mac OSX. For Windows users, the individual components of this protocol, such as
112 BLAST, software used for RNA-Seq analysis, programming language R and Python, all have Windows
113 compatible executable files and can be used under Windows environments. In this protocol, we will
114 install NCBI BLAST for the homology search step (Section 2.2), STAR for read mapping and
115 featureCounts for counting reads (Section 2.6), and the R programming language and several packages for
116 RNA-Seq and comparative transcriptome analysis (Section 2.7).

117

118 **2.1 Set up folder structure for data analysis.**

119

120 To facilitate reproducible and effective computational analysis [33,34], we suggest that the user create a
121 folder structure (**Figure 1**) such that the raw data, processed data, results, and scripts for data processing
122 can be organized into their respective folders. In this protocol, the reader can use the following commands
123 to create the recommended folder structure.

124

```
125     $ cd ATH_GMA  
126     $ mkdir raw_data processed_data scripts results software  
127     $ mkdir processed_data/bam processed_data/rc
```

128

129 Sequence and annotation files from databases should be downloaded to the “**raw_data**” folder. Software
130 tools that will be used in this analysis can be saved and installed in the “**software**” folder. We recommend

131 the reader to create a folder named “**bin**” under the **software** folder such that the executable files can be
132 copied to “**software/bin**” folder and add “**software/bin**” to the PATH environmental variable under the
133 Linux environment. For experienced Linux users, software can also be installed in a user specified folder
134 such as ~/bin or in a system wide folder. The reader can download scripts in github into the “**scripts**”
135 folder. Intermediate output will be generated in the “**processed_data**” folder, and major input and output
136 files for visualization will be saved in the “**results**” folder.

137
138 All scripts for this step are provided in “Section2.1_setup_directory.sh” in the “scripts” folder. The reader
139 can set up the folder structure (**Figure 1**) using the following command.

```
140  
141     $ cd ATH_GMA  
142     $ sh ./scripts/Section2.1_setup_directory.sh
```

143
144 **[Figure 1 near here]**

145 146 **2.2 Software installation**

147 We provide a script to download and install tools for RNA-seq analysis; readers can run the script in the
148 project folder.

```
149  
150     $ cd ATH_GMA  
151     $ sh ./scripts/Section2.2_download_softwares.sh
```

152
153 A successfully installed tool will return version information when it is run only with a “-v” or a “—
154 version” option.

155

156 **Install NCBI BLAST for identification of homologous genes.** BLAST is a sequence similarity search
157 tool [35]. The latest version of NCBI BLAST can be downloaded from the NCBI ftp site using the
158 following link: <ftp://ftp.ncbi.nlm.nih.gov/blast/executables/LATEST/>. This folder contains precompiled
159 executable files and installation files for Windows, Mac OSX, and Linux platforms. Because finding
160 orthologous genes at a genome scale is computationally intensive, it is recommended to use a Linux
161 workstation or computing cluster to perform the BLAST analysis.

162

163 For Linux users, the current pre-compiled executable is `ncbi-blast-2.6.0+-x64-linux.tar.gz`.

164 For Mac users, the current installation file is `ncbi-blast-2.6.0+.dmg`.

165 For Windows users, the current installation file is `ncbi-blast-2.6.0+-win64.exe`.

166

167 A later version of BLAST should work as well with minor changes in the command line options. For
168 Windows and Mac users, double click the downloaded file to install the program. For Linux users, one
169 can use “`tar -xvf ncbi-blast-2.6.0+-x64-linux.tar.gz`” to extract the archive file. After extracting the files,
170 move the executable files to a folder in the Linux search path.

171

172 **Install tools for RNA-Seq data download.** The following shows a sample script to download sra-tools
173 and fastq-dump to download the raw sequencing data. The sequence read archive (SRA) database
174 provides sra-toolkit, which is a suite of easy to use computational tools to download data from the
175 database. To download the raw data from the SRA database, one needs to first install the sra-toolkit and
176 use the fastq-dump utility program based on the SRA ids.

177

178 `$ cd ATH_GMA/software`

179 `$ wget http://ftp-trace.ncbi.nlm.nih.gov/sra/sdk/current/sratoolkit.current-centos_linux64.tar.gz`

180 `$ tar -xzf sratoolkit.current-centos_linux64.tar.gz`

181 `$./sratoolkit.2.8.2-1-centos_linux64/bin/fastq-dump --version`

182

183 **Install tools for RNA-Seq data analysis.** We will install the STAR [36] and featureCounts [37] software
184 tools. STAR is a read mapper, and featureCounts can count the number of reads mapped to each gene in
185 the genome. Both software tools were used here due to their speed and accuracy [38,39]. Other alternative
186 mappers can be used, and there are excellent review papers [39–41] that compare and summarize these
187 different bioinformatics tools.

188

189 To download and install STAR and featureCounts, run the following scripts in the project folder.

190

```
191 $ cd Proj_CompTS_ATH_GMA/software
```

```
192 $ wget https://github.com/alexdobin/STAR/archive/2.5.2b.tar.gz
```

```
193 $ tar -xzf 2.5.2b.tar.gz
```

```
194 $ STAR-2.5.2b/bin/Linux_x86_64_static/STAR --version
```

```
195 $ wget https://sourceforge.net/projects/subread/files/subread-1.5.1/subread-1.5.1-Linux-
```

```
196 x86\_64.tar.gz/download
```

```
197 $ tar -zxvf download
```

```
198 $ subread-1.5.1-Linux-x86_64/bin/featureCounts -v
```

199

200 **2.3 Install R, DESeq2, and edgeR packages for RNA-Seq data analysis.**

201 R is a programming language and environment for statistical data analysis [42]. We will use R to
202 summarize RNA-Seq reads and to generate FPKM data. To install R, the reader should go to the
203 Comprehensive R Archive Network (CRAN) (<https://cran.r-project.org>) to download the installer packages
204 for their Windows, Mac OSX, or Linux system. For Linux users, R can be installed using the command
205 line, and platform dependent package management systems. For example, to install R in CentOS 7 Linux,
206 the user should simply type:

207

208 \$ sudo yum install R

209

210 Scripts for installing R packages are provided in:

211

212 Section2.3_install_r_packages.R

213

214 To install DESeq2 and edgeR, the user should follow the instructions for these respective packages. These
215 two packages are part of the Bioconductor repository such that the installation should be performed using
216 the Bioconductor installation script. The following commands are executed under the R environment and
217 these commands are preceded by ">". For commands that are executed under Linux terminals, these
218 commands are preceded by "\$".

219

220 > source('https://bioconductor.org/biocLite.R')

221 > biocLite('DESeq2')

222 > biocLite('edgeR')

223

224 The installation script will detect the dependency of these two packages and install other required
225 packages accordingly.

226

227 To install the OrthoClust package, the user should download the script for the OrthoClust package.

228

229 > setwd("./software")

230 > install.packages("OrthoClust_1.0.tar.gz", repos=NULL, type="source")

231

232 **2.4 Download protein and genome sequences for Arabidopsis and soybean.**

233 Sample scripts for download are provided in “Section2.4_download_data.sh”. All protein-coding
234 sequences and genomic sequences for Arabidopsis can be downloaded from the Araport web site
235 (www.araport.org). Araport is a data portal for Arabidopsis genomic research that hosts the latest
236 genomic sequences and genome annotations for this model organism [43]. The web site requires free
237 registration to access the download link to the protein sequences and genome annotation files. As of July
238 2017, the current version of the protein sequences file is “Araport11_genes.201606.pep.fasta.gz”. This
239 name will likely be different for future versions of the protein sequences. We recommend that users
240 download the latest version of the protein sequences, and record the actual download date and version of
241 the sequence files for the purpose of reproducibility. The latest version of the genome sequence of
242 Arabidopsis is “TAIR10_Chr.all.fasta.gz”. This file is unlikely to change because the genome assembly of
243 Arabidopsis is likely to remain the same in the future. The latest version of the gene annotation file is
244 “Araport11_GFF3_genes_transposons.201606.gtf.gz”.

245
246 All protein-coding sequences for soybeans can be downloaded from the DOE phytozome database
247 (https://phytozome.jgi.doe.gov/pz/portal.html#!bulk?org=Org_Gmax). Phytozome is a data portal for
248 plant and microbial genomes that hosts dozens of sequenced plant genomes and gene annotations [44].
249 This web site also requires free registration before data downloading. The latest version of soybean
250 protein sequences is version 2.0 (downloaded in July 2017). The protein sequences and genomic
251 sequences are “Gmax_275_Wm82.a2.v1.protein.fa.gz” and “Gmax_275_v2.0.fa.gz”. These names are
252 likely to change with future versions of the genome and proteome annotation. The latest version of the
253 gene annotation file is “Gmax_275_Wm82.a2.v1.gene_exons.gff3.gz”.

254
255 These files are in compressed fasta format and require de-compression before use. Under the Linux
256 command line, the following command can be used to de-compress these “*.gz”.

257
258 `$ gunzip Araport11_genes.201606.pep.fasta.gz`

259 \$ gunzip Gmax_275_Wm82.a2.v1.protein.fa.gz

260

261 **2.5 Download raw data from published RNA-Seq experiments**

262 Raw sequencing data can be downloaded from the NCBI Sequence Read Archive (SRA)

263 (<https://www.ncbi.nlm.nih.gov/sra>). The embryo developmental data sets for Arabidopsis and soybean

264 can be found in two bioprojects (PRJNA301162 for Arabidopsis and PRJNA197379 for soybean). For

265 the Arabidopsis samples, RNA-Seq data were collected in triplicates at seven time points (7, 8, 10, 12, 13,

266 15, and 17 days after pollination). For the soybean samples, RNA-Seq data were collected in triplicates at

267 ten time points (5, 10, 15, 20, 25, 30, 35, 40, 45, and 55 days, day 0 of the time course is 12 to 17 days

268 after anthesis). Each sample is represented by a unique GSM id; for example, the three replicates of 7

269 days old Arabidopsis embryo samples are GSM1930276, GSM1930277, and GSM1930278. All 41

270 samples from this experiment are stored under a unique GSE id, GSE74692. Each sample is also

271 represented by a unique SRA id. For example, the three replicates of 7 days old Arabidopsis embryo

272 samples are SRR2927328, SRR2927329, and SRR2927330 from PRJNA301162.

273

274 \$ fastq-dump --split-3 SRR2927328 --outdir ./raw_data

275

276 We suggest that the reader download the data into the raw data folder for further processing. To download

277 large numbers of data sets, prepare a text file with all SRR ids for one species and run the following script

278 in the project folder.

279

280 \$ cd ATH_GMA

281 \$ sh ./scripts/Section2.5_download_fastq.sh ./raw_data/PRJNA301162.txt ATH

282 \$ sh ./scripts/Section2.5_download_fastq.sh ./raw_data/PRJNA197379.txt GMA

283

284 Depending on the size of sequencing data and network speed, this step may take a few hours. We provide
285 a test file “PRJNAtest.txt” for the user to test the execution time for downloading one file. The time for
286 downloading the entire data set can be estimated based on downloading this single file. We also provide
287 the FPKM data for this particular data set so that the users do not need to download the original data to
288 perform the analysis in this protocol. To perform the analysis using provided FPKM file, the user can
289 start the analysis from Section 3.4.

290

291 **3. Methods**

292 **3.1 Comparative transcriptome analysis overview.**

293 This protocol provides details of comparative transcriptome analysis between two species. We not only
294 compute sequence similarity between protein coding genes in two species, we also integrate the gene
295 expression patterns of these genes from two different species under similar biological processes. There
296 are three major steps in this analysis (**Figure 2**): 1) identify homologous genes between two species; 2)
297 generate a gene expression data matrix and a co-expression network in each species; 3) perform cross
298 species comparisons of gene homology and expression patterns. For each of these steps, multiple
299 bioinformatics tools are available. This protocol will provide a basic workflow for each of the steps and
300 the reader can substitute individual steps with other tools (**See Note 4.3**).

301

302 **[Figure 2 near here]**

303

304 **3.2 Identifying homologous genes between species.**

305 **3.2.1 Identification of homologous pairs using BLAST.**

306 Analysis in this section can be performed using the following command:

307

```
308 $ cd ATH_GMA
```

```
309 $ sh ./scripts/Section3.2.1_BLAST.sh
```

310

311 **Step 1.** Merge the Arabidopsis protein fasta file and soybean protein fasta file using this Linux command:

312

```
313 $ cat Araport11.pep.fasta GLYMA2.pep.fasta > ATHGMA.pep.fasta
```

314

315 **Step 2.** Create the BLAST database:

316

```
317 $ makeblastdb -in ATHGMA.pep.fasta \
```

```
318     -out ATHGMA.blastdb \
```

```
319     -dbtype prot \
```

```
320     -logfile makeblastdb.log
```

321

322 The option “-in” specifies the input file name of the merged protein fasta file. “-out” specifies the

323 BLAST database file name. “-dbtype” indicates the database is a protein database. “-logfile” is for

324 recording error messages in case the process fails.

325

326 **Step 3.** Perform the BLAST search.

327 The Linux command used in this step is:

328

```
329 $ blastp -evalue 0.00001 \
```

```
330     -outfmt 6 -db ATHGMAX.blastdb \
```

```
331     -query ATHGMA.fasta > ATHGMA.pep.blastout
```

332

333 The option “-evalue” specifies the E value threshold. “-outfmt” is set to be 6, which is tab delimited

334 format. “-db” is set to be the BLAST database built in step 3. “-query” uses the merged protein fasta files

335 as input. The results of BLAST analysis are written in a file named ATHGMAX.pep.blastout.

336

337 The output includes the following 12 tab-separated columns “**qseqid sseqid pident length mismatch**
338 **gapopen qstart qend sstart send evalue bitscore**”. The meaning of these columns can be found using
339 the BLAST help manual. The columns that will be used in downstream analysis are **qseqid** (query
340 sequence id), **sseqid** (subject sequence id), and **evalue** (E value). We will filter BLAST results and only
341 keep homologous genes with BLAST E value $< 1e-5$ [3,26].

342

343 **3.2.2 Obtaining reciprocal best hit (RBH) genes**

344

345 Reciprocal best BLAST hit (RBH) and its variants are commonly used methods to identify homologous
346 genes in two species [45–49]. To identify RBH genes between any two species, the BLAST results from
347 protein sequence alignment were first parsed to identify the best BLAST hit for each soybean protein in
348 the Arabidopsis protein lists. For each soybean protein, there is at most one best BLAST hit protein in the
349 Arabidopsis proteome. For each of the Arabidopsis proteins identified in the first step, the best BLAST
350 hit of each protein in the soybean proteome is also identified. If this best hit is also the original
351 homologous gene found in the first step, this pair of proteins is defined to constitute an RBH pair.

352

353 For genes with multiple isoforms and potentially multiple protein sequences, we performed the BLAST
354 analysis at the isoform level and then collapsed all the isoforms for each gene to find the best match. In
355 fact, a large fraction of the isoforms in both Arabidopsis and soybean do not change their protein coding
356 sequences, the difference being found in the UTR regions of the transcripts being compared. This is
357 consistent with published results in Arabidopsis and soybean [50,51]. We developed a Python script that
358 can identify RBH genes from the above two species from BLAST results. The user can download this
359 script from the github repository. To perform the analysis the user can use the following commands:

360

361 `$ cd ATH_GMA`

362 \$ sh ./scripts/Section3.2.2_RBH.sh

363

364 Although RBH genes are widely used in comparative genomic analysis, other methods can be used to
365 identify homologous genes for downstream analysis (see **Note 4.3**). An example file
366 (ARATH2GLYMA.RBH.subset.txt) of RBH genes is provided. The user can use this file to perform the
367 following analysis without running the RBH script.

368

369 **[Table 1 near here]**

370

371 **3.3 Gene expression data processing.**

372 Gene expression quantification includes three main steps: 1) read mapping; 2) read counting and 3)
373 FPKM calculation. For this analysis, we follow a published protocol for expression processing [50].

374

375 **Step 1.** Create genome index by STAR.

376 RNA-Seq reads have to be mapped to the respective reference genomes. To use STAR to map reads to the
377 reference genome, the user needs to build a genome index using the following commands.

378

379 \$ cd ATH_GMA

380 \$ sh ./scripts/Section3.3.Step1.MakeIndex.sh

381

382 The following commands are used to create a genome index for Arabidopsis.

383

384 \$ WORKDIR=\$(pwd)

385 \$ IDX=\$WORKDIR/raw_data/ATH_STAR-2.5.2b_index

386 \$ GNM=\$WORKDIR/raw_data/TAIR10_Chr.all.fasta

387 \$ GTF=\$WORKDIR/raw_data/Araport11_GFF3_genes_transposons.201606.gtf

388 \$ STAR --runMode genomeGenerate \

```
389         --genomeDir $IDX \  
390         --genomeFastaFiles $GNM \  
391         --sjdbGTFfile $GTF
```

392
393 The option “--runMode” indicates that the command is to create a genomic index. “--genomeDir”
394 specifies the file name for the genome index. “--genomeFastaFiles” indicates the input fasta file for
395 genomic sequences. “--sjdbGTFfile” is to provide a genome annotation file when creating the genomic
396 index. A genome index will be created for each species.

397
398 **Step 2.** Read mapping by STAR.

399 After creating genome indexes, the user needs to use STAR to map reads from each sample to the
400 reference genome to generate a read mapping file using the following commands.

```
401  
402     $ cd ATH_GMA  
403     $ sh ./scripts/Section3.3.Step2.Mapping.ATH.sh  
404     $ sh ./scripts/Section3.3.Step2.Mapping.GMA.sh
```

405
406 The “Section3.3.Step2.Mapping.ATH.sh” is to map all Arabidopsis reads. The
407 “Section3.3.Step2.Mapping.GMA.sh” is to map all Soybean reads. In the SRA database, each sample has a
408 unique SRR id. The following commands show one example of such SRR ids (SRR2927328).
409 SRR2927328_1 and SRR2927328_2 represent two ends of paired reads.

```
410  
411     $ STAR --genomeDir $IDX \  
412           --readFilesIn $WORKDIR/raw_data/SRR2927328_1.fastq.gz  
413           $WORKDIR/raw_data/SRR2927328_2.fastq.gz \  
414           --outFileNamePrefix $WORKDIR/processed_data/bam/SRR2927328/SRR2927328 \  
415           --outSAMtype BAM SortedByCoordinate
```


416

417 The option “--genomeDir” specifies the file name for the genome index. “--readFilesIn” indicates the
418 input fastq files for RNA-seq reads. Two files are provided for paired-end reads. “--outFileNamePrefix” is
419 to provide the directory for output data. “--outSAMtype BAM” indicate the output file should be a bam
420 file. “SortedByCoordinate” set the output data to be sorted by the order of where the read is mapped to the
421 chromosome.

422

423 **Step 3.** Read counting with featureCounts.

424

425 To count reads with featureCounts, the user can use the following command:

426

```
427     $ cd ATH_GMA
```

```
428     $ sh ./scripts/Section3.3.Step3.ReadCount.ATH.sh
```

```
429     $ sh ./scripts/Section3.3.Step3.ReadCount.GMA.sh
```

430

431 For this step, featureCounts will calculate how many reads map to each gene region. For simplicity, we
432 only count uniquely mapped reads and only summarize read counts at the gene level. Other software can
433 be used to summarize expression at isoforms levels. The following commands are for counting reads for a
434 single file.

435

```
436     $ WORKDIR=$(pwd)
```

```
437     $ GTF=$WORKDIR/raw_data/Araport11_GFF3_genes_transposons.201606.gtf
```

```
438     $ BAM=$WORKDIR/processed_data/bam
```

```
439     $ RC=$WORKDIR/processed_data/rc
```

```
440     $ featureCounts -t exon \
```

```
441         -g gene_id \
```

```
442         -p \
```

```
443         -a $GTF \  
444         -o $RC/SRR2927328.readcount.txt \  
445         $BAM/SRR2927328/SRR2927328Aligned.sortedByCoord.out.bam
```

446
447 The option “-t exon” indicates that only reads mapped to exons are counted. The option “-p” indicate the
448 input reads are paired-end reads. The option “-a” provides the location of the genome annotation file. The
449 option “-o” specifies the output file location. The last parameter is the file name of the read mapping file
450 (bam file).

451
452 **Step 4.** FPKM calculation using DESeq2 and edgeR.

453
454 For this step, R scripts will be used to summarize gene expression level in fragments per kilo-basepairs
455 per million reads (FPKM). To calculate FPKM, we performed the following five steps: 1) merging read
456 counts from different files into one single file; 2) differential expression analysis using DESeq2; 3) data
457 normalization. 4) FPKM calculation and 5) average FPKM calculation across replicates. These steps can
458 be performed using a unified sh (shell) script: NGS_RNA-seq_CalcFPKM.R, which is provided in the
459 github repository of this project. To run this script, the user needs to provide a table that summarizes the
460 replicate structure of the samples. Example tables (PRJNA301162.csv for Arabidopsis and
461 PRJNA197379.csv for soybean) are provided in the “processed_data” folder.

462
463 To run the unified R script for FPKM calculation, use the following commands:

```
464  
465     $ cd ATH_GMA  
466     $ Rscript ./scripts/Section3.3.Step4.FPKM.R ./processed_data/fpkm/GMA  
467     $ Rscript ./scripts/Section3.3.Step4.FPKM.R ./processed_data/fpkm/ATH
```

468

469 This script requires multiple input files to be present in the working directory. These files include a file
470 that describes the design matrix of the experiment and the read count files generated in Step 3. More
471 descriptions of the input file formats are included in the annotation of the R script.

472

473 **Step 5. Co-expression Networks from gene expression profiles**

474 Expression data will be summarized and converted to gene co-expression networks. The input data
475 include data matrices with averaged and normalized FPKM values. In this protocol, we use genes in
476 metabolic pathways that are essential to seed development. Other methods can be used to filter genes
477 before the analysis, for example, only keep genes with high variations across conditions. Finally, gene co-
478 expression matrices were calculated for each species. We use the cut-off with p value < 0.001 and
479 Pearson Correlation Coefficient > 0.99 to generate co-expression networks. To generate co-expression
480 networks, the following commands were used.

481

```
482 $ cd ATH_GMA
```

```
483 $ Rscript ./scripts/Section3.3.Step5_FPKM2NETWORK.R
```

484

485 **3.4. Identify orthologous co-expressed clusters using OrthoClust**

486 **3.4.1 Overview of the OrthoClust method.**

487 Simple approaches can be used to identify conserved co-expression genes across different species. For
488 example, one can first cluster gene expression in two species separately, and, for each pair of cluster
489 combinations, one can find whether the pairs of clusters share significantly large numbers of homologous
490 genes using appropriate statistical tests such as Fisher's exact test. OrthoClust [29] is a global approach
491 in which the process of co-expression clustering finding and homology detection is integrated into the
492 same objective function. The objective function H is defined as

493

$$H = - \left(\sum_{i,j \in S_1} \Lambda_{ij}^1 \delta_{\sigma_i \sigma_j} + \sum_{i,j \in S_2} \Lambda_{ij}^2 \delta_{\sigma_i \sigma_j} + \kappa \sum_{(i,j') \in O(S_1, S_2)} w_{ij'} \delta_{\sigma_i \sigma_{j'}} \right)$$

494
495
496 where S_N is the sets of genes for a species and a subscript of S ($N = 1$ or 2) corresponds to the species
497 respectively. i and j are individual genes of a species or nodes on a network. Λ_{ij}^N denotes a modularity
498 score from gene i and j , that is a difference between the real number of edges and the expected number of
499 edges. $\delta_{\sigma_i \sigma_j}$ is for a module label. If i and j have the same module label, $\delta_{\sigma_i \sigma_j} = 1$, and, if not, $\delta_{\sigma_i \sigma_j} = 0$. A
500 coupling constant, κ controls overall impact of orthology relations on the objective function, and a weight,
501 $w_{ij'}$ is for orthology relations coming from the number of orthologous genes between two species. The
502 objective function H will return lower values when orthologous genes are assigned into the same module.
503
504 This approach translates orthologous co-expression finding into a network module finding problem. The
505 objective function includes three components: two components represent the goodness of the expression
506 clustering results and one component represents the effect of homologous genes across species. The
507 parameter κ can be adjusted to increase or decrease the contribution of homologous genes in the
508 clustering processes. The effects of using different co-expression thresholds and parameter κ are
509 discussed in Note 4.4.

511 **3.4.2 Steps for OrthoClust analysis.**

512
513 To perform OrthoClust analysis, we require three input data files: 1) the gene co-expression network from
514 soybean; 2) the gene co-expression network from Arabidopsis; and 3) the orthologous gene pairs between
515 two species.

516

517 These files require a specific format for the OrthoClust engine to analyze. The user can use the following
518 R command to perform the clustering analysis

```
519  
520 > library(OrthoClust)  
521 > OrthoClust2(Eg1=GMX_edgelist, Eg2=ATH_edgelist, \  
522 list_orthologs=GA_orthologs, kappa=3)
```

523
524 We provide a wrapper script that will read three input files: a list of edges from Arabidopsis, a list of
525 edges from soybean, and a list of RBH gene pairs from two species. To perform OrthoClust analysis, the
526 user can simple use the following commands:

```
527  
528 $ cd ATH_GMA  
529 $ Rscript ./scripts/Section3.4.Step1_OrthoClust.R
```

530
531 This script will generate three files. “Orthoclust_Results.csv” contains information regarding modules
532 assignment for each gene. “Orthoclust_Results_Summary.csv” includes number of genes assigned to
533 each module. “Orthoclust_Results.RData” contains multiple R objects that will be used in the
534 visualization step.

535
536 **[Table 3 near here]**

537 **3.4.3 Visualization of OrthoClust results as a network.**

538
539 To visualize OrthoClust results, we use Cytoscape, a network visualization platform to analyze biological
540 networks and to integrate multiple data into networks such as gene expression profiles or annotation [52].
541 We used module 8 from the previous step as an example. There are three input files: 1) soybean co-
542 expression network edge list for genes in module 8, 2) Arabidopsis co-expression network edge list for

543 genes in module 8, and 3) RBH list for genes in module 8. To generate these files for Cytoscape
544 visualization, the user can use the following command.

545

```
546 $ cd ATH_GMA
```

```
547 $ Rscript ./scripts/Section3.4.Step2_CytoscapeInput.R
```

548

549 **[Figure 3 near here]**

550

551 **Step 1.** To Import three files on Network Browser, we can first start from the Cytoscape menu bar “File” >
552 “import” > “Network” > “File”. After you select one of three input files, the popup window with “Import
553 Network From Table” title appears. You can see two columns with gene names in the middle of the
554 window. Next, to change attributes of columns, click the first line of each column and choose either
555 “Source Node” or “Target Node” from the menu. Since three edge lists do not have direction, the two
556 columns from each input file can be assigned into either source or target nodes. After that, we change an
557 option for column names from “Advanced Options” at the bottom left of the window. On the new popup
558 window, we can uncheck “Use first line as column names”, since we do not have headers in the input files.
559 Finally, you can see two column names, “Column1” and “Column2” with different icons of attributes,
560 and the remaining parts of the preview are gene names. You can repeat these steps for each of the input
561 files.

562

563 **Step 2.** With three imported networks, we can integrate data sets of co-expression networks with
564 homologous relations using the Union function. To do that, select three network on the network tab on the
565 control panel (click one network and click the other two networks while pressing Command), and move
566 to Cytoscape’s menu bar “Tools” > “Merge” > “Networks”.

567

568 In the popup window for “Advanced Network Merge”, we should choose the “Union” button, select three
569 networks from “Available Networks”, and then click the right-facing arrow acting for “Add Selected”.
570 After that you can find that three networks are now on “Networks to Merge”, and you can click “Merge”
571 button to merge three networks.

572
573 The name of the merged network will appear with the total number of merged nodes and edges on the
574 Network tab on the control, and usually it is automatically visualized on the Cytoscape canvas.

575
576 **Step 3.** To express properties of networks (species information, source of edges such as co-expression
577 networks or homologous relations), we can customize visual attributes of the merged network. To do that,
578 on the Select tab on the control panel, we can click the “+” icon below the “Default filter” and choose
579 “Column Filter” to add the new condition. From the “Choose column” drop-down list, you can select
580 “Node: name” or “Edge: name” and type a prefix of each species (“AT” for Arabidopsis genes, or
581 “Glyma” for soybean genes). This filter applies to visualization of the merged automatically, so you can
582 see highlighted nodes on the Cytoscape canvas.

583
584 There are several ways to change visualization properties of the selected components. First, we can set
585 “Bypass Style” for the selected nodes or edges such as “Fill Color” and “Size” for properties of nodes, or
586 “Stroke Color” and “Line Type” for properties of edges. To do this, move your mouse pointer on one of
587 the highlighted nodes, right-click, and then select “Edit” > “Bypass Style” > “Set Bypass to Selected
588 Nodes” on the popup menu. The control panel on the left side will be automatically changed to the “Style”
589 tab, and you can see three subtabs: “Node”, “Edge”, and “Network” on the bottom of the interface.
590 Second, we can apply different Layouts with these selected nodes or all nodes from Cytoscape menu bar
591 “Layout”.

592

593 As an example of the network with module 8, nodes and edges from soybean and Arabidopsis genes were
594 switched to green and orange colors respectively. To highlight genes of interest, we used thicker double
595 lines for edges and blue color for nodes. We separated genes into four groups according to their input files
596 and species (Arabidopsis genes from RBH results or not, and soybean genes from RBH results or not),
597 and layout each of them with Degree Sorted Circle Layout (Figure 3).

598

599 **3.4.4 Visualization of OrthoClust results as expression profiles.**

600 We also provide scripts to directly visualize gene expression patterns for orthologous co-expression
601 modules (**Figure 4**). This figure is generated by the script “Section3.4.Step2_CytoscapeInput.R”. In this
602 module, most soybean genes are tightly clustered. Some Arabidopsis genes are tightly clustered (close to
603 the black line) whereas other Arabidopsis genes are not. This result shows that many genes in the soybean
604 co-expression cluster change their expression patterns in Arabidopsis, suggesting potential functional
605 divergence of these genes. In contrast, many genes that are RBH pairs in the two species have similar
606 expression patterns. For example, one gene (AT5G52560, green line) that is related to the raffinose
607 biosynthetic pathway has a similar decreasing expression pattern as its RBH gene (Glyma.04G245100) in
608 soybean.

609 **[Figure 4 near here]**

610

611 **4. Notes**

612 **4.1 Software installation:**

613 The git software is installed in most Linux systems by default. If git is not installed in your system, please
614 refer to <https://git-scm.com> for installation instructions.

615

616 **4.2 Code blocks.** All code blocks started with “\$” are command line scripts that should be executed under
617 a Linux terminal. All code blocks started with “>” are command line scripts that should be executed
618 under an interactive R programming language console.

619

620 **4.3** For each of these steps, multiple bioinformatics tools are available. This protocol will provide a basic
621 workflow for each of the steps and the reader can substitute individual steps with other tools. For example,
622 in searching for homologous genes, several other alternative tools such as OMA or OrthoFinder [10,11]
623 can be used instead of BLAST. A comprehensive comparison of these tools is out of the scope of this
624 chapter. Some databases or tools provide pre-computed homologous genes [8,12]. Additional steps must
625 be performed to ensure that the gene ids from OMA, OrthoFinder, or PLAZA match the gene ids used in
626 the expression analysis.

627

628 **4.4** Many genes in both species were not included in the RBH gene lists. This is because the criterion for
629 identifying RBH genes is highly stringent, as it requires that both genes in two species be the best BLAST
630 hit in their respective species. This can be relaxed to identify k-best-hits in two species [6]. We have
631 developed a script that can generate k-best-hits using BLAST results between any two species
632 (`OrthologousGenes_OneWayTopNBestHit.py`).

633

634 **4.5** Effect of different parameters in OrthoClust analysis. We analyzed how different parameters affect the
635 results of this analysis. We focus on two major parameters (**Figure 5**): the Pearson Correlation
636 Coefficient (PCC) threshold that was used to convert co-expression data to networks, and the kappa
637 parameter that was used in OrthoClust analysis.

638

639 **[Figure 5 near here]**

640

641 The kappa parameter is used to adjust the relative importance of the co-expression edges and homologous
642 edges in network module finding algorithms. When kappa equals zero, the module finding method only
643 finds co-expression modules and does not consider the effects of homologous edges. When kappa is set to
644 be higher than zero, homologous edges will be included in the module finding objective function. This

645 can be verified by comparing the numbers of modules found by $\kappa = 0$ to numbers of modules found
646 by $\kappa > 0$. The numbers of modules found by $\kappa = 1$ is 2 to 3 times the numbers of modules found
647 by $\kappa = 0$. This result suggests that including homologous edges generates more modules across
648 species, because, when $\kappa = 0$, all modules are from the same species. Comparing the numbers of
649 modules from $\kappa = 2$ with $\kappa = 1$, and $\kappa = 3$ with $\kappa = 2$ suggest that increasing κ can
650 further increase the number of modules.

651

652 The PCC threshold also affects the number of modules identified. For the same κ value, a higher
653 PCC threshold always leads to more modules. This is expected as a co-expression network with higher
654 PCC threshold contains fewer edges. Because of the reduced number of edges, the network is less
655 connected and can be break into more modules as compared to the network generated with lower PCC
656 threshold.

657

658 **Ethics approval and consent to participate**

659 Not applicable

660

661 **Consent for publication**

662 Not applicable

663

664 **Availability of data and materials**

665 The datasets and software supporting the conclusions of this article are available in the Github repository
666 (<https://github.com/LiLabAtVT/CompareTranscriptome>).

667

668 **Competing interests**

669 The authors declare no competing interests.

670

671 **Authors' contributions**

672 JL and SL designed the analysis. RG provided the original data and interpreted the biological results. LH

673 edited the manuscript and provided suggestions to improve the methods. JL developed the methods. SL

674 and JL wrote the manuscript.

675

676 **Acknowledgments and Funding**

677 This work is partly supported by Virginia Soybean Board.

678

679 **Literature Cited**

680 1. Movahedi S, Van Bel M, Heyndrickx KS, Vandepoele K. Comparative co-expression analysis in plant

681 biology. *Plant. Cell Environ.* 2012;35:1787–98.

682 2. Van Bel M, Proost S, Wischnitzki E, Movahedi S, Scheerlinck C, Van de Peer Y, et al. Dissecting

683 Plant Genomes with the PLAZA Comparative Genomics Platform. *Plant Physiol.* 2012;158:590–600.

684 3. Movahedi S, Van de Peer Y, Vandepoele K. Comparative Network Analysis Reveals That Tissue

685 Specificity and Gene Function Are Important Factors Influencing the Mode of Expression Evolution in

686 *Arabidopsis* and Rice. *Plant Physiol.* 2011;156:1316–30.

687 4. Prince SJ, Joshi T, Mutava RN, Syed N, Joao Vitor M dos S, Patil G, et al. Comparative analysis of the

688 drought-responsive transcriptome in soybean lines contrasting for canopy wilting. *Plant Sci. Elsevier*

689 Ireland Ltd; 2015;240:65–78.

690 5. Ware D, Jaiswal P, Ni J, Pan X, Chang K, Clark K, et al. Gramene: a resource for comparative grass

691 genomics. *Nucleic Acids Res.* 2002;30:103–5.

692 6. Mutwil M, Klie S, Tohge T, Giorgi FM, Wilkins O, Campbell MM, et al. PlaNet: combined sequence

693 and expression comparisons across plant networks derived from seven species. *Plant Cell.* 2011;23:895–

694 910.

- 695 7. Ruprecht C, Mendrinna A, Tohge T, Sampathkumar A, Klie S, Fernie AR, et al. FamNet: A framework
696 to identify multiplied modules driving pathway diversification in plants. *Plant Physiol.*
697 2016;170:pp.01281.2015.
- 698 8. Li LL, Jr CJS, Roos DS. OrthoMCL: Identification of Ortholog Groups for Eukaryotic Genomes.
699 *Genome Res.* 2003;13:2178–89.
- 700 9. Jensen LJ, Julien P, Kuhn M, von Mering C, Muller J, Doerks T, et al. eggNOG: automated
701 construction and annotation of orthologous groups of genes. *Nucleic Acids Res.* 2007;36:D250–4.
- 702 10. Altenhoff AM, Kunca N, Glover N, Train C-M, Sueki A, Pili ota I, et al. The OMA orthology
703 database in 2015: function predictions, better plant support, synteny view and other improvements.
704 *Nucleic Acids Res.* 2015;43:D240–9.
- 705 11. Emms DM, Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons
706 dramatically improves orthogroup inference accuracy. *Genome Biol.* 2015;16:157.
- 707 12. Proost S, Bel M Van, Vanechoutte D, Peer Y Van De, Mueller-roeber B, Vandepoele K. PLAZA 3 .
708 0□: an access point for plant comparative genomics *Dirk Inz e.* 2015;43:974–81.
- 709 13. Proost S, Fostier J, Witte D De, Dhoedt B, Demeester P, Peer Y Van De, et al. genomic homology in
710 extremely large data sets. 2012;40.
- 711 14. Lyons E, Pedersen B, Kane J, Alam M, Ming R, Tang H, et al. Finding and Comparing Syntenic
712 Regions among Arabidopsis and the Outgroups Papaya, Poplar, and Grape: CoGe with Rosids. *PLANT*
713 *Physiol.* 2008;148:1772–81.
- 714 15. Berri S, Abbruscato P, Faivre-Rampan O, Brasileiro ACM, Fumasoni I, Satoh K, et al.
715 Characterization of WRKY co-regulatory networks in rice and Arabidopsis. *BMC Plant Biol.* 2009;9:1–
716 22.
- 717 16. Wang Y, Feng L, Zhu Y, Li Y, Yan H, Xiang Y. Comparative genomic analysis of the WRKY III
718 gene family in populus, grape, arabidopsis and rice. *Biol. Direct.* 2015;10:48.
- 719 17. Yao X, Ma H, Wang J, Zhang D. Genome-Wide Comparative Analysis and Expression Pattern of
720 TCP Gene Families in Arabidopsis thaliana and Oryza sativa. *J. Integr. Plant Biol.* 2007;49:885–97.

- 721 18. Netotea S, Sundell D, Street NR, Hvidsten TR. ComPIEx: conservation and divergence of co-
722 expression networks in *A. thaliana*, *Populus* and *O. sativa*. *BMC Genomics*. 2014;15:106.
- 723 19. Roulin A, Auer PL, Libault M, Schlueter J, Farmer A, May G, et al. The fate of duplicated genes in a
724 polyploid plant genome. *Plant J*. 2013;73:143–53.
- 725 20. Provart NJ, Alonso J, Assmann SM, Bergmann D, Brady SM, Brkljacic J, et al. 50 years of
726 Arabidopsis research: highlights and future directions. *New Phytol*. 2016;209:921–44.
- 727 21. Horan K, Jang C, Bailey-Serres J, Mittler R, Shelton C, Harper JF, et al. Annotating Genes of Known
728 and Unknown Function by Large-Scale Coexpression Analysis. *Plant Physiol*. 2008;147:41–57.
- 729 22. Schneider A, Aghamirzaie D, Elmarakeby H, Poudel AN, Koo AJ, Heath LS, et al. Potential targets of
730 VIVIPAROUS1/ABI3-LIKE1 (VAL1) repression in developing *Arabidopsis thaliana* embryos. *Plant J*.
731 2016;85:305–19.
- 732 23. Aghamirzaie D, Nabiyouni M, Fang Y, Klumas C, Heath L, Grene R, et al. Changes in RNA Splicing
733 in Developing Soybean (*Glycine max*) Embryos. *Biology (Basel)*. 2013;2:1311–37.
- 734 24. Kilian J, Whitehead D, Horak J, Wanke D, Weinl S, Batistic O, et al. The AtGenExpress global stress
735 expression data set: protocols, evaluation and model data analysis of UV-B light, drought and cold stress
736 responses. *Plant J*. 2007;50:347–63.
- 737 25. Wang L, Czedik-Eysenberg A, Mertz RA, Si Y, Tohge T, Nunes-Nesi A, et al. Comparative analyses
738 of C₃ and C₄ photosynthesis in developing leaves of maize and rice. *Nat. Biotechnol*. Nature Publishing
739 Group; 2014;32:1158–65.
- 740 26. Patel R V, Nahal HK, Breit R, Provart NJ. BAR expressolog identification: expression profile
741 similarity ranking of homologous genes in plant species. *Plant J*. 2012;71:1038–50.
- 742 27. Junker A, Hartmann A, Schreiber F, Bäumlein H. An engineer's view on regulation of seed
743 development. *Trends Plant Sci*. 2010;15:303–7.
- 744 28. Chaudhury AM, Koltunow A, Payne T, Luo M, Tucker MR, Dennis ES, et al. Control of Early Seed
745 Development. *Annu. Rev. Cell Dev. Biol*. 2001;17:677–99.
- 746 29. Yan K-K, Wang D, Rozowsky J, Zheng H, Cheng C, Gerstein M. OrthoClust: an orthology-based

- 747 network framework for clustering data across multiple species. *Genome Biol.* 2014;15:R100.
- 748 30. Palla G, et al. Directed network modules. *New J. Phys.* 2007;9:186.
- 749 31. Malliaros FD, Vazirgiannis M. Clustering and Community Detection in Directed Networks: A Survey.
- 750 *Phys. Rep.* 2013;533:86.
- 751 32. Smoot ME, Ono K, Ruscheinski J, Wang P-L, Ideker T. Cytoscape 2.8: new features for data
- 752 integration and network visualization. *Bioinformatics.* 2011;27:431–2.
- 753 33. Sandve GK, Nekrutenko A, Taylor J, Hovig E. Ten Simple Rules for Reproducible Computational
- 754 Research. Bourne PE, editor. *PLoS Comput. Biol.* 2013;9:e1003285.
- 755 34. Osborne JM, Bernabeu MO, Bruna M, Calderhead B, Cooper J, Dalchau N, et al. Ten Simple Rules
- 756 for Effective Computational Research. Bourne PE, editor. *PLoS Comput. Biol.* 2014;10:e1003506.
- 757 35. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J. Mol.*
- 758 *Biol.* 1990;215:403–10.
- 759 36. Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, et al. STAR: ultrafast universal
- 760 RNA-seq aligner. *Bioinformatics.* 2013;29:15–21.
- 761 37. Liao Y, Smyth GK, Shi W. FeatureCounts: An efficient general purpose program for assigning
- 762 sequence reads to genomic features. *Bioinformatics.* 2014;30:923–30.
- 763 38. Su Z, Łabaj PP, Li S, Thierry-Mieg J, Thierry-Mieg D, Shi W, et al. A comprehensive assessment of
- 764 RNA-seq accuracy, reproducibility and information content by the Sequencing Quality Control
- 765 Consortium. *Nat. Biotechnol.* 2014;32:903–14.
- 766 39. Conesa A, Madrigal P, Tarazona S, Gomez-Cabrero D, Cervera A, McPherson A, et al. A survey of
- 767 best practices for RNA-seq data analysis. *Genome Biol.* 2016;17:13.
- 768 40. Dillies MA, Rau A, Aubert J, Hennequet-Antier C, Jeanmougin M, Servant N, et al. A comprehensive
- 769 evaluation of normalization methods for Illumina high-throughput RNA sequencing data analysis. *Brief.*
- 770 *Bioinform.* 2013;14:671–83.
- 771 41. Steijger T, Abril JF, Engström PG, Kokocinski F, The RGASP Consortium, Hubbard TJ, et al.
- 772 Assessment of transcript reconstruction methods for RNA-seq. *Nat. Methods.* 2013;10:1177–84.

- 773 42. R Core Team. R: A Language and Environment for Statistical Computing. 2017;
- 774 43. Krishnakumar V, Hanlon MR, Contrino S, Ferlanti ES, Karamycheva S, Kim M, et al. Araport: the
775 Arabidopsis information portal. *Nucleic Acids Res.* 2015;43:D1003-9.
- 776 44. Goodstein DM, Shu S, Howson R, Neupane R, Hayes RD, Fazo J, et al. Phytozome: a comparative
777 platform for green plant genomics. *Nucleic Acids Res.* 2012;40:D1178-86.
- 778 45. Ward N, Moreno-Hagelsieb G. Quickly Finding Orthologs as Reciprocal Best Hits with BLAT, LAST,
779 and UBLAST: How Much Do We Miss? de Crécy-Lagard V, editor. *PLoS One.* 2014;9:e101850.
- 780 46. Fulton DL, Li YY, Laird MR, Horsman BGS, Roche FM, Brinkman FSL. Improving the specificity of
781 high-throughput ortholog prediction. *BMC Bioinformatics.* 2006;7:270.
- 782 47. Dessimoz C, Boeckmann B, Roth ACJ, Gonnet GH. Detecting non-orthology in the COGs database
783 and other approaches grouping orthologs using genome-specific best hits. *Nucleic Acids Res.*
784 2006;34:3309–16.
- 785 48. Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin E V, et al. The COG database:
786 an updated version includes eukaryotes. *BMC Bioinformatics.* 2003;4:41.
- 787 49. O'Brien KP, Remm M, Sonnhammer ELL. Inparanoid: a comprehensive database of eukaryotic
788 orthologs. *Nucleic Acids Res.* 2005;33:D476-80.
- 789 50. Li S, Yamada M, Han X, Ohler U, Benfey PN. High resolution expression map of the Arabidopsis
790 root reveals alternative splicing and lincRNA regulation. *Dev. Cell.* 2016;in press:508–22.
- 791 51. Aghamirzaie D, Collakova E, Li S, Grene R. CoSpliceNet: a framework for co-splicing network
792 inference from transcriptomics data. *BMC Genomics.* 2016;17:845.
- 793 52. Cline MS, Smoot M, Cerami E, Kuchinsky A, Landys N, Workman C, et al. Integration of biological
794 networks and gene expression data using Cytoscape. *Nat. Protoc.* 2007;2:2366–82.
- 795
- 796 **Key Reference**
- 797 Yan et al., 2014. See above.

798 A methodology to cluster integrated data from co-expression profile for each species and from
799 homologous relationships between multiple species.

800

801 **Internet Resources**

802 <https://www.araport.org>

803 *Arabidopsis information portal*

804 https://phytozome.jgi.doe.gov/pz/portal.html#!bulk?org=Org_Gmax

805 *Genomics resource page of Glycine max Wm82.a2.v1 in Phytozome*

806 <https://git-scm.com>

807 *Git software Home page*

808 <https://github.com/LiLabAtVT/CompareTranscriptome.git>

809 *Github page for this tutorial*

810 <https://www.ncbi.nlm.nih.gov/sra>

811 *NCBI Sequence Read Archive (SRA) Home page*

812 <https://github.com/alexdobin/STAR>

813 *Github page of STAR*

814 <http://bioinf.wehi.edu.au/subread-package/>

815 *The Subread package Web page*

816 <https://cran.r-project.org>

817 *The Comprehensive R Archive Network Web page*

818

819 **Figure Captions**

820 **Figure 1.** Folder structure for data analysis.

821 **Figure 2.** A workflow of comparative transcriptome analysis between soybean and Arabidopsis. It is
822 composed of three major parts: identification of ortholous pairs between two species using BLAST,
823 RNA-seq analysis to get co-expression networks, and running OrthoClust to cluster genes with
824 orthologous relations. Blue fonts indicates softwares or scripts used in this workflow.

825 **Figure 3.** Visualization of module 8 from OrthoClust result. In this network, Circle 1 and 4 stand for
826 groups of genes from Arabidopsis and soybeans that do not have orthology in the other species and only
827 co-expression partner from the same species. Circle 2 and 3 denote genes have orthologous partner in the
828 other species as well as their co-expression partners from the same species. Green nodes are genes from
829 Arabidopsis, and red from soybean. Edges from co-expression network of Arabidopsis are green, and
830 those of soybeans are red. Black double lined edges indicate homologous pairs between soybean and
831 Arabidopsis genes. Four genes from raffinose biosynthesis pathways are highlighted in blue color and
832 their homologous pairs have thicker edges.

833 **Figure 4.** Expression plots of genes from Arabidopsis and soybean bellowing to one of modules of
834 OrthoClust result. One example of homologous genes in Arabidopsis and soybeans are AT5G52560 and
835 Glyma.04G245100 are highlighted in green.

836 **Figure 5.** Effect of different correlation cutoff and κ values on the number of modules in orthoClust
837 analysis.

838

839 **Table Captions**

840 **Table 1.** Results of Identified Orthologous Genes.

841

842 **Table 2.** Examples of input data files for OrthoClust analysis. There are three inputs: two co-expression
 843 networks of (A) soybean and (B) Arabidopsis, (C) orthologous pairs between soybean and Arabidopsis.

844 **Table 3.** Top 10 OrthoClust results sorted by the total number of genes from a module. OrthoClust was
 845 performed with parameters $\kappa=3$, gene co-expression correlation cutoff \geq correla and homologous pairs
 846 obtained from RBH Blast.

847

848 **Tables**

849 **Table 1.** Results of Identified Orthologous Genes.

Species	Soybean	Arabidopsis
Number of proteins	48,375	24,148
(Total number of gene models)	(56,044)	(37,336)
Blast results in each species	1,086,080	1,081,623
(Query: Blast DB)	(Soybean: Arabidopsis)	(Arabidopsis: Soybean)
Number of RBH genes in each species	13,024	13,024
Number of 5 best hit in each species	208,343	112,819

850

851 **Table 1.** Examples of input data files for OrthoClust analysis. There are three inputs: two co-expression
 852 networks of (A) soybean and (B) Arabidopsis, (C) orthologous pairs between soybean and Arabidopsis.

(A)		(B)		(C)	
row	column	row	column	Soybean gene	Arabidopsis gene
Glyma.01G006400	Glyma.01G016500	AT1G01540	AT1G05350	Glyma.01G001300	AT2G07050
Glyma.01G021300	Glyma.01G021400	AT1G06040	AT1G06150	Glyma.01G005800	AT4G29310
Glyma.01G019400	Glyma.01G022500	AT1G01720	AT1G07400	Glyma.01G006100	AT4G26300
Glyma.01G015400	Glyma.01G026700	AT1G05230	AT1G07570	Glyma.01G010100	AT1G32090
Glyma.01G025100	Glyma.01G026700	AT1G02660	AT1G08230	Glyma.01G015400	AT2G35470

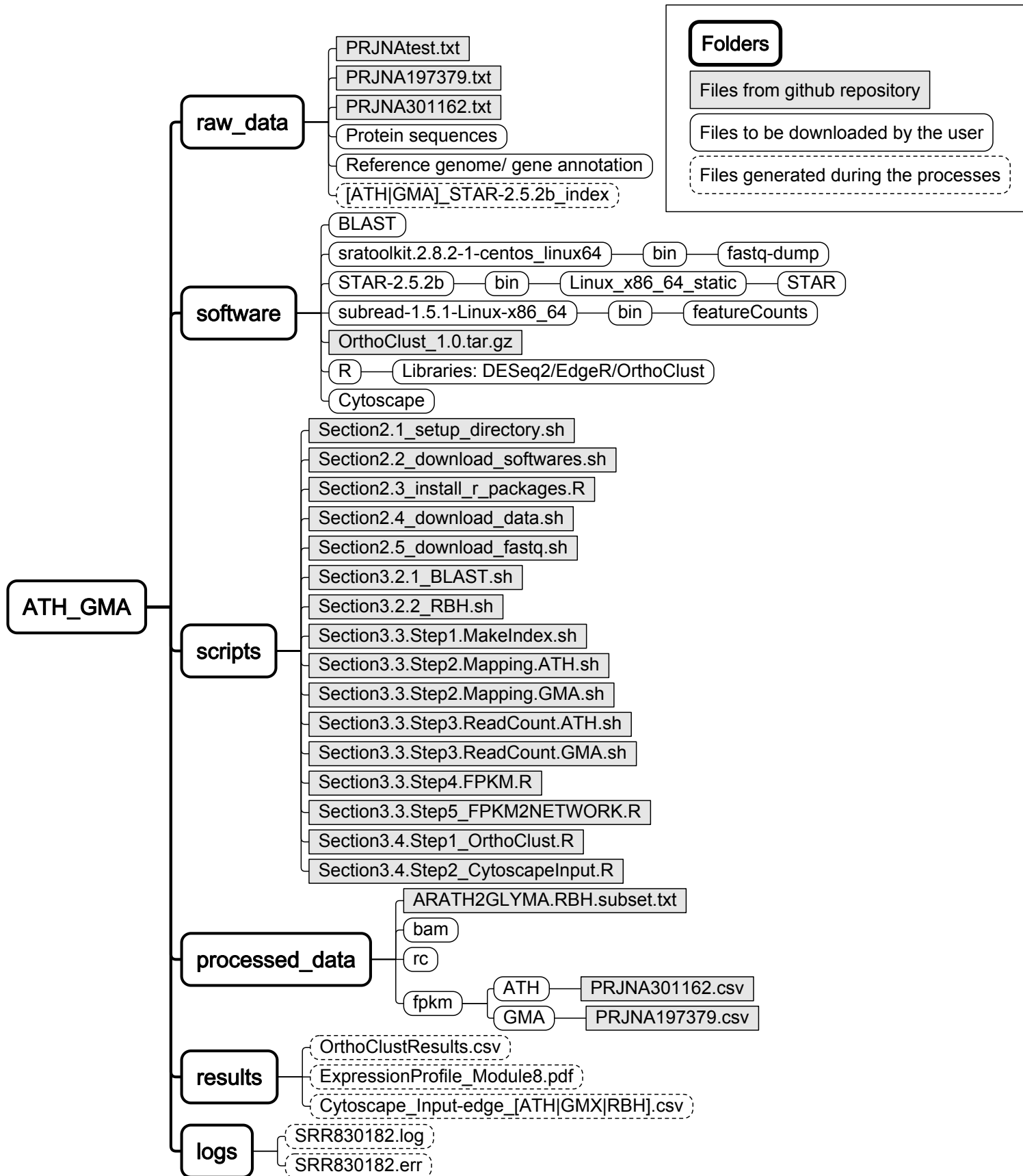
Glyma.01G025100 Glyma.01G028900 AT1G01090 AT1G08510 Glyma.01G019400 AT5G65670

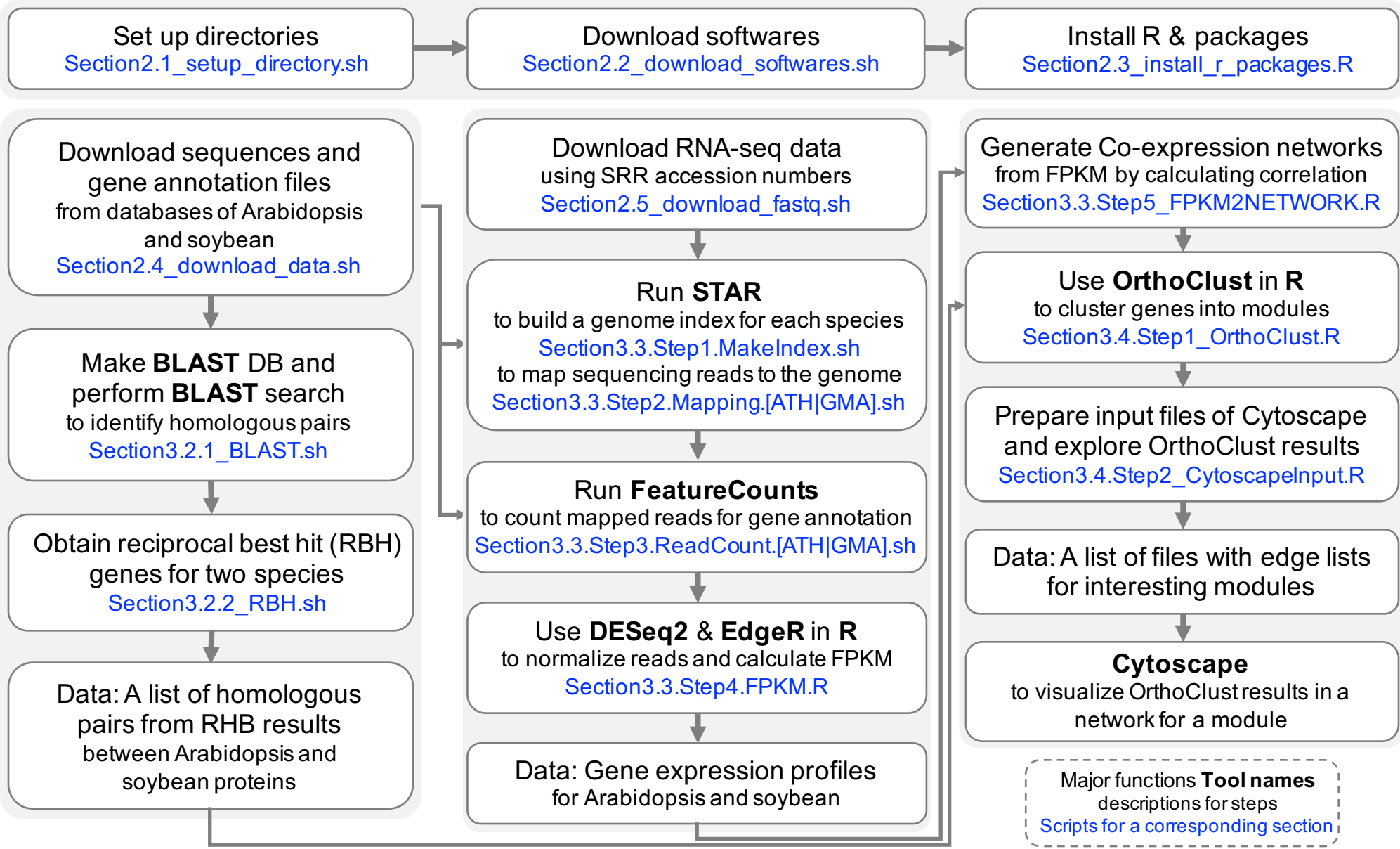
853

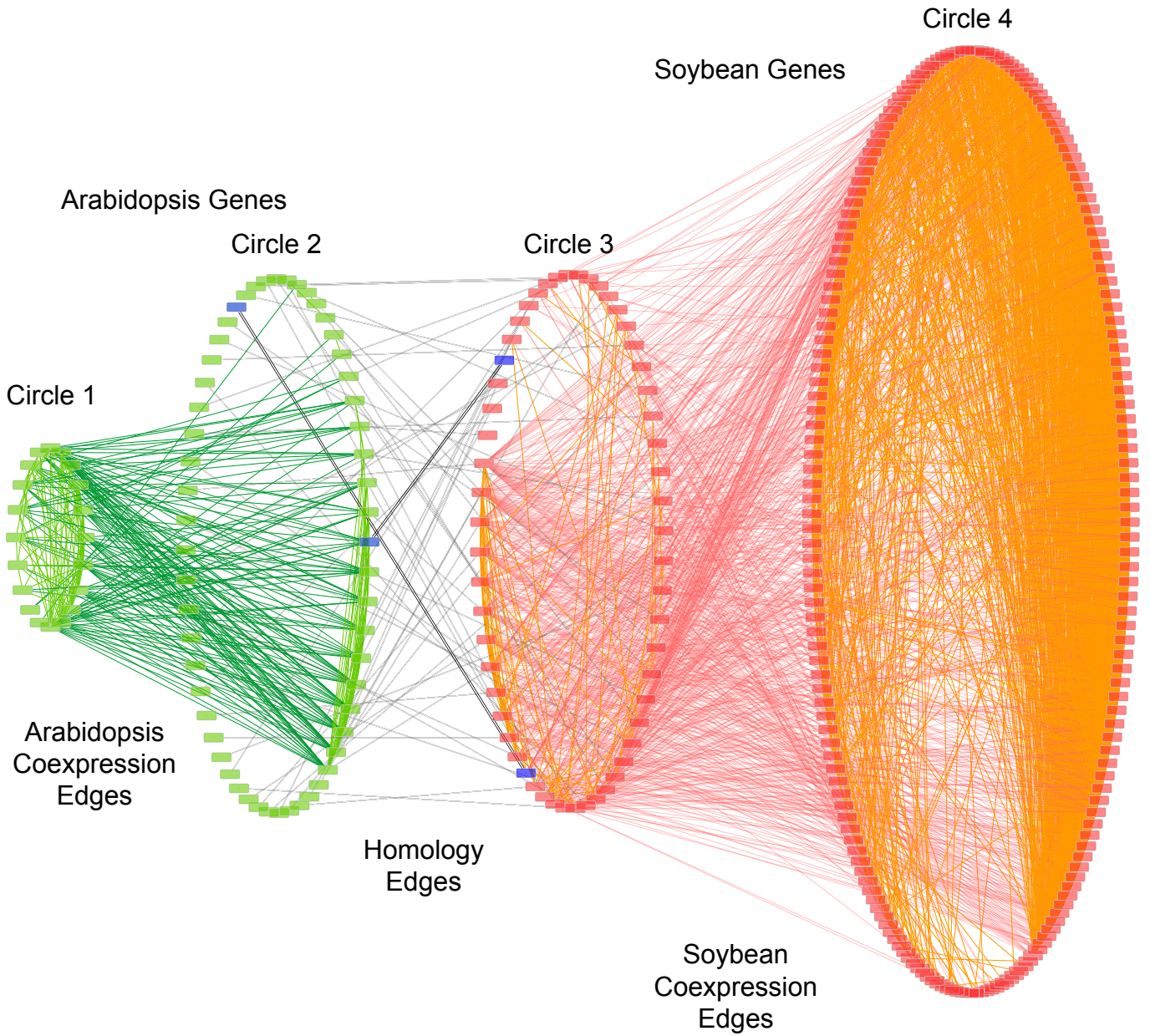
854 **Table 3.** Top 10 OrthoClust results sorted by the total number of genes from a module. OrthoClust was
 855 performed with parameters $\kappa=3$, gene co-expression correlation cutoff \geq correla and homologous pairs
 856 obtained from RBH Blast.

No.	Module ID	Total number of genes from a module	The number of genes from soybean	The number of genes from Arabidopsis
1	2	352	273 (77.6%)	79 (22.4%)
2	8	331	255 (77.0%)	76 (23.0%)
3	39	297	174 (58.6%)	123 (41.4%)
4	1	253	214 (84.6%)	39 (15.4%)
5	3	245	207 (84.5%)	38 (15.5%)
6	187	215	56 (26.0%)	159 (74.0%)
7	224	212	53 (25.0%)	159 (75.0%)
8	57	192	110 (57.3%)	82 (42.7%)
9	113	147	38 (25.9%)	109 (74.1%)
10	19	45	39 (86.7%)	6 (13.3%)

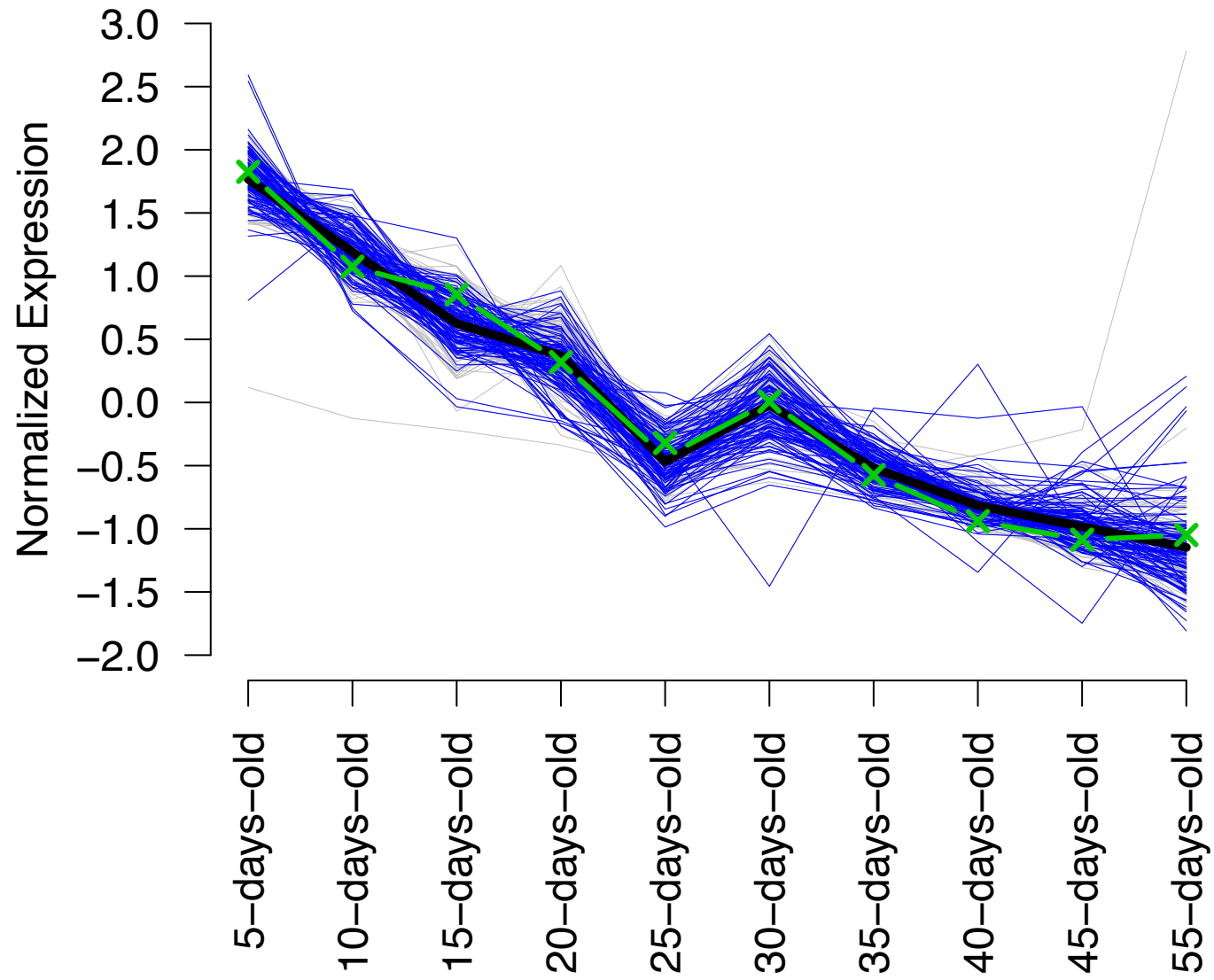
857



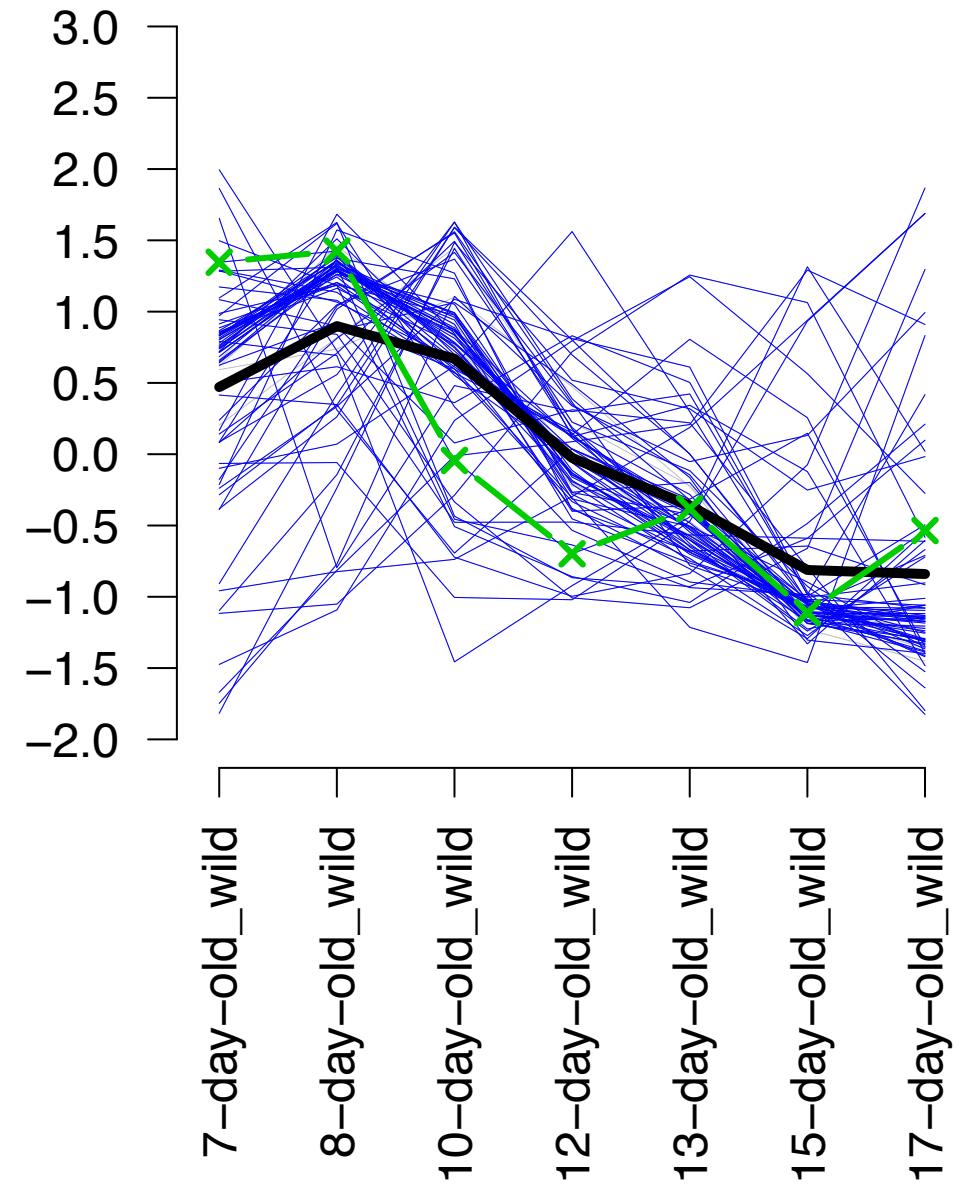




Soybean



Arabidopsis



— Cluster Average Expression —x— Example homologous genes — Homogous Genes — Non-homogous Genes

The number of modules

