

Selection of suitable reference genes for qRT-PCR studies during SE initial dedifferentiation in cotton of different SE capability

Cao Ai Ping¹, Shao Dong Nan¹, Cui Bai Ming¹, Zheng Yin Ying¹, Sun jie^{Corresp. 2}

¹ Shihezi University, Colleges of Life Science, Shihezi, Shihezi, China

² Shihezi University, The Key Laboratory of Oasis Eco-Agriculture, Shihezi, Shihezi, China

Corresponding Author: Sun jie
Email address: sunjie@shzu.edu.cn

Analysis of gene expression level by RNA sequencing (RNA-seq) has a wide range of biological purposes in various species. Real-time fluorescent quantitative PCR (qRT-PCR) evaluated gene expression levels and validated transcriptomic, which will depend on the stably expressed reference genes for normalization of the gene expression level under specific situations. In this study, 15 candidate genes were selected from transcriptome datasets during somatic embryogenesis (SE) initial dedifferentiation in *Gossypium hirsutum* L. of different SE capability. To evaluate the stability of those genes, geNorm, NormFinder and BestKeeper were used. The results revealed that *ENDO4* and *18srRNA* could be as appropriate reference genes under all conditions. The stability and reliability of the reference genes were further tested through comparison of qRT-PCR results and RNA-seq data, as well as evaluation of the expression profiles of auxin-responsive protein (*AUX22*) and ethylene-responsive transcription factor (*ERF17*). In summary, the results of our study indicate the most suitable reference genes for qRT-PCR during three induction stages in four cotton species.

Selection of suitable reference genes for qRT-PCR studies during SE initial dedifferentiation in cotton of different SE capability

Cao Aiping¹, Shao Dongnan¹, Cui Baiming¹, Zheng Yinying¹, Sun Jie^{2*}

¹ Colleges of Life Science, Shihezi University, Shihezi, China

² The Key Laboratory of Oasis Eco-Agriculture, Shihezi University, Shihezi, China

*Corresponding: sunjiexj@shzu.edu.cn

Analysis of gene expression level by RNA sequencing (RNA-seq) has a wide range of biological purposes in various species. Real-time fluorescent quantitative PCR (qRT-PCR) evaluated gene expression levels and validated transcriptomic, which will depend on the stably expressed reference genes for normalization of the gene expression level under specific situations. In this study, 15 candidate genes were selected from transcriptome datasets during somatic embryogenesis (SE) initial dedifferentiation in *Gossypium hirsutum* L. of different SE capability. To evaluate the stability of those genes, geNorm, NormFinder and BestKeeper were used. The results revealed that *ENDO4* and *18srRNA* could be as appropriate reference genes under all conditions. The stability and reliability of the reference genes were further tested through comparison of qRT-PCR results and RNA-seq data, as well as evaluation of the expression profiles of auxin-responsive protein (*AUX22*) and ethylene-responsive transcription factor (*ERF17*). In summary, the results of our study indicate the most suitable reference genes for qRT-PCR during three induction stages in four cotton species.

Keywords: qRT-PCR; *Gossypium hirsutum* L.; reference gene; gene expression analysis; RNA sequencing

Introduction

Gene expression level analysis is crucial in many fields of biological research ¹. In particular, RNA sequencing (RNA-seq) has become the prevalent method for gene transcriptomes expression analysis of various species ^{2,3}. Moreover, real-time fluorescent quantitative PCR (qRT-PCR) has also been widely applied to evaluate gene expression levels and to validated transcriptomic ^{4,5}, and qRT-PCR provides several advantages, including higher sensitivity, reproducibility and specificity ⁶. However, qRT-PCR is affected by inaccurate quantification of RNA, and RNA quality contributes to non-specific variations ⁷. Thus, to avoid bias, it is crucial to select stably expressed reference genes for normalization of the gene expression level under certain conditions. Reliable reference genes are usually selected from housekeeping genes because of their stable expression ¹, as well as their involvement in basic cellular function and cell processes, such as cell structure formation, cytoskeletal protein

formation, and ribosomal subunit synthesis ^{8,9}.

Recent reports have demonstrated that some reference genes fluctuated under different conditions in different species, and that they do not have stable expression. For example, the halophyte *Halostachys caspica* under salt and drought stress, *EF1α* and *TUB3* would be suitable for normalizing gene expression data, respectively, while *UBC10* was most stable under both stresses combined ⁹. In *Caragana intermedia*, *UNK2*, *PP2A* and *SAND* were appropriated for the most suitable reference genes across all tested conditions, while *UNK2*, *SAND* and *EF1α* were the most stable reference genes for salt-treated leaves, and *UNK2* and *SAND* provide superior transcript normalization for salt-treated roots. Additionally, *TIP41* and *PP2A* were the most suitable for PEG-treated (osmotic) leaves, while *UNK1*, *UNK2* and *PP2A* were stably expressed for PEG-treated roots, and *SAND* and *EF1α* exhibited the most perfect expression in cold-treated leaves. *SAND*, *TIP41* and *PP2A* should be sufficient to normalize in heat-treated leaves¹⁰. In *Gracilaria lemaneiformis* under temperature stress, *GAPDH*, *ITS2*, *CR*, and *18SrRNA* were the optimal reference genes for various treatments applied at 8°C, while *eIF* and *ACT* were optimal for 32°C, and *GAPDH*, *EF1α* and *ACT* were ideal for different temperature treatments ¹⁰. In *Hedera helix* L., *40S* was the optimal reference genes with ABA treatment, while under different tissues and various cold stress conditions, *HhSE* and *Hhb-AS* were perfected for reference genes ¹¹.

Cotton is a major raw material for the textile industry and also a source of oil. SE is an important process in cotton molecular breeding with *Agrobacterium*-mediated transformation. However, only a few cotton varieties have been successfully used in genetic engineering in vitro via SE. ¹². What more, the majority of the cotton cultivars in the production are poor regeneration via SE ¹³, and many factors can effect SE, including culture conditions and tissue background; therefore, most studies conducted to date have focused on identifying molecular mechanisms and genes critical for SE ^{14,15}. However, no systematic analyses of reference genes can be used in SE development.

In this study, we selected the following 15 candidate genes from transcriptome datasets of *Gossypium hirsutum* L based on a q-value ≥ 0.05 , a small coefficient of variation (CV) of $FPKM \geq 2$ and $|\log_2 \text{FoldChange}| < 1$: *18srRNA*, *ARF1* (ADP-ribosylation factor 1), *ARF2* (ADP-ribosylation factor 2), *EF1α* (elongation factor 1-alpha), *ENDO4* (endonuclease 4), *ERF3A* (eukaryotic peptide chain release factor 3A), *IF4E2* (eukaryotic translation initiation factor isoform 4E-2), *NUB1* (NEDD8 ultimate buster 1), *PTBP3* (polypyrimidine tract-binding protein homolog 3), *RPAB5* (DNA-directed RNA polymerases I, II, and III subunit), *T2FB* (transcription initiation factor IIF, beta subunit), *TAF11* (TBP-associated factor 11), *UBE4* (U-box domain-containing protein), *UBC7* (ubiquitin carrier protein 7) and *UFD1* (ubiquitin fusion degradation 1). Three common software programs, geNorm ¹⁶, NormFinder ¹⁷ and BestKeeper ¹⁸, were served to calculate the expression stability of candidate reference genes. Moreover, *ERF17* and *AUX22* were the main genes used to determine the stability and reliability of the

71 reference gene.

72 Materials and Methods

73 Plant materials

74 Four cultivars of *Gossypium hirsutum* L, YZ1, R15, X33 and X42, were investigated, among
75 which YZ1 and R15 have a relatively high SE differentiation rate and therefore comprise the
76 main transgenic material¹⁹. Although X33 and X42 are the major commercial cultivars in
77 Xinjiang, China, they have a low rate of differentiation during SE compared with YZ1 and R15²⁰.
78 Specimens of all cultivars have been conserved in our laboratory.

79 Total RNA extraction and cDNA reverse transcription

80 Total RNA was extracted and cDNA was synthesized from each sample using the
81 Purelink™ RNA Mini Kit (Life Technologies, Carlsbad, CA, USA) and the PrimerScript RT
82 reagent Kit (Perfect Real Time) (RR037Q, Takara) following the manufacturer's protocols. The
83 detailed steps were established following our previously described method²⁰.

84 Selection of reference genes and design of primer

85 We performed high throughput RNA-Seq on four varieties of *Gossypium hirsutum* L
86 cultivars at three treatment stages using the Illumina Hiseq™ 2000 platform. To evaluate gene
87 expression stability, $q\text{-value} \geq 0.05$, $\text{FPKM} \geq 2$ and $|\log_2\text{FoldChange}| < 1$ were used as the criteria
88 for selection of candidate housekeeping genes at all sampling points²¹. We selected 15
89 candidate housekeeping genes from the cotton transcriptome. Specific primers were designed
90 based on sequences of the 15 reference genes using NCBI/Primer-BLAST ([https://](https://www.ncbi.nlm.nih.gov/tools/primer-blast/)
91 www.ncbi.nlm.nih.gov/tools/primer-blast/) according to the following parameters: PCR
92 product size of 100–150 bp, primer melting temperatures (T_m) of 57°C–60°C and PCR primer
93 pairs separated by at least one intron on the corresponding genomic DNA. All primer pairs were
94 synthesized by BGITECH. The PCR products were tested by electrophoreses on 1% agarose
95 gels.

96 qRT-PCR analysis

97 qRT-PCR was executed in 96-well plates in a LightCycler® 480 Real-Time PCR System
98 (Roche) with a SYBR Green-based PCR assay. Reactions with a total volume of 10 μL

including 1 μ L of template (first-strand cDNA) , 0.4 μ L each of 10 μ M forward and reverse gene-specific primers , 5 μ L of 2 \times SYBR Premix Ex Taq II (TLi RaseH Plus) (Takara, Dalian, China) and 3.2 μ L of ddH₂O. The qRT-PCR conditions were as follows: initial denaturation at 95°C for 30s, followed by 45 cycles of 95°C for 10s, 60°C for 10 s and 72°C for 10s. qRT-PCR analysis was tested in three biological replicates. Additionally, three technical replicates were used for each qRT-PCR analysis, and the correlation coefficient (R^2) and amplification efficiency (E) were counted by a standard curves with the diluted series on the basis of the diluted cDNA series²². The PCR efficiency was detected by the equation ($E = (10^{[-1/\text{slope}]} - 1) \times 100\%$)¹⁶.

Data analysis

The raw Cq values are listed in the Supplementary Table S1. Three common software programs, geNorm¹⁶, NormFinder¹⁷ and BestKeeper¹⁸, were applied to calculate the expression stability of candidate reference genes. For geNorm and NormFinder, Cq values were converted into relative quantities according to the formula: $2^{-\Delta C_t}$ (ΔC_t = the corresponding Cq value – minimum Cq)²³. The BestKeeper calculations are directly based on raw Cq values of each gene. Microsoft Excel was used for variance analyses.

Validation of Reference Genes

To test the stability and reliability of the reference genes, we compared the qRT-PCR results and RNA-seq data. Moreover, the expression profiles of *ERF17* and *AUX22* were analyzed and compared with the FPKM values obtained in the RNA-seq data.

Results

Isolation of cotton references genes

Fifteen candidate reference genes were screened from transcriptome datasets of *Gossypium hirsutum* L., which were obtained using the Illumina HiseqTM 2000 platform. q-value \geq 0.05, FPKM \geq 2, |log₂FoldChange|<1 and relatively lower CV of FPKM were set as the higher criteria at all sampling points for selection of candidate housekeeping genes. Fig. 1 shows a heat map of candidate reference genes at each pairwise comparison.

Verification of primer specificity and PCR amplification efficiency

The target genes were cloned with cDNA as the template using specific primers, respectively. Gel electrophoresis confirmed that all specific primers amplified a PCR product, and each specific primer pairs was validated by the presence of a single peak in the melting curve. The PCR amplification efficiency and correlation coefficients (R^2) of the 15 reference genes were counted according to the slopes of the standard curves in four cotton species. As shown in Table 1, the qRT-PCR amplification efficiency ranged from 96.13% for TAF11 to 120.30% for UBE4, while the R^2 value ranged from 0.9677 to 1.0977.

Candidate gene expression profile analyses at different treatment stages between cotton species

The expression pattern of 15 reference genes was identified for four cotton species in three treatment stages by qRT-PCR. The expression levels of the reference genes were calculated by the number of cycles (C_q), which referred to the amplification related fluorescence to reach a specific threshold level of test, with a smaller C_q indicating higher gene expression. The 15 reference genes displayed a relatively wide range of mean C_q values of 17.36 to 22.61 (Fig.2, Table 2). *ERF3A* showed the highest expression level at the X42 of 3d; however, *END04* expressed the lowest level at R15 of 0h. Similarly, *18srRNA* showed the least gene expression variation, while *UBE4* showed the greatest variation among samples.

Reference gene expression stability analyses

To identify the most stably expressed gene for cotton qRT-PCR normalization, a gene expression stability study was conducted using three publicly available statistical tools, geNorm, NormFinder and BestKeeper.

a) geNorm analysis

geNorm can be used to check the most suitable number of reference genes. All candidate reference genes were estimated based on their M values¹⁶. The value of M for each gene was analyzed stability of reference genes, where $M=1.5$ is the standard of stability for reference gene expression, and $M<1.5$ is the expression level of reference genes thought to be stable. Thus, a lower M value represents higher reference gene expression stability¹⁶. As shown in

Table 3 and Fig. 3, the analysis results showed through geNorm implied that *18s rRNA*, *END04* and *ARF1* were the most suit genes for four different cotton species in three treatment stages, while *UBC7* (0.38), *UBC7* and *UBE4* were the least consistently expressed. For the cotton species of YZ1, *ARF1* had the lowest M value ($M=0.24$), whereas *TAF11* displayed the highest M value ($M=0.46$) under the same conditions. All genes were individually analyzed under R15, and the most stable genes were *END04*, *IF4E2* and *18s rRNA*. Among X33, the most stable genes were *END04*, *ARF2* and *ARF1*. The most stable genes in X42 were *ARF2*, *T2FB* and *UFD1*.

The pairwise variation V value ($V_{n/n+1}$) of one gene with the others can also establish the optimal number of reference genes for normalization. In general, $V_{n/n+1}$ is the cut-off value. If $V_{n/n+1} = 0.15$, the best reference genes number for correct normalization should be optimal $n+1$; if $V_{n/n+1} < 0.15$, the reference genes number should be n ²⁴. As shown in Fig. 3, the $V_{2/3}$ value was below 1.5 in the current study, the result suggested that two reference genes were required for normalization.

168 b) NormFinder analysis

The NormFinder software is a Visual Basic application tool, which be used to establish the expression stabilities of reference genes based on the stability value (Sv). NormFinder analysis has some varies between geNorm analysis, which takes into account intra- and intergroup variations for normalization factor (NF) calculations¹⁷. Genes with lower average expression stability values have higher stability, which is thought to indicate a stable reference gene. As shown in Table 4, the results of NormFinder analysis were relatively consistent with those of geNorm. The least stable genes were as follows: *18s rRNA*, *END04* and *ARF1* in total, three induction of four cotton cultivars, *PRAB5*, *ARF1* and *18s rRNA* for YZ1; *END04*, *IF4E2* and *18s rRNA* for R15; *END04*, *ARF2* and *18s rRNA* for X33; and *ARF2*, *T2FB* and *UFD1* for X42. The least stable gene was *TaF11*, which was consistently found in all groups.

179 c) Bestkeeper analysis

The Bestkeeper software¹⁸ can be used to analyze the stability and expression of reference genes according to the coefficient of variance (CV), standard deviation (SD) and correlation coefficient (R). Reference genes are considered to be stable when they have a high R value and low SD and CV values. Additionally, if $SD > 1$, the gene was considered unacceptable⁹. As shown in Table 5, *YZ1*, *END04*, *18s rRNA* and *ARF1* were relatively stable expression genes under all conditions. *TAF11* showed a low $CV \pm SD$ value and a lower R value, indicating it was an unstable gene. For *R15*, *18srRNA*, *UBC7* and *ARF1* were relatively suitable as reference genes. Overall, *END04*, *18srRNA* and *PTBP3* were the most

stable, with high R values and low CV± SD values in X33. Among the cotton of X44, *T2FB*, *ARF2* and *EDD04* were considered acceptable reference genes because of the strong correlation and low CV± SD value.

Validation of Reference Gene

To test the stability and reliability of the reference genes, we compared qRT-PCR results and RNA sequencing (RNA-seq) data. Additionally, two target genes, *AUX22* and *ARF17*, were selected to further validate the reference genes. The relative expression levels of *AUX22* and *ARF17* were calculated using the validated reference genes and compared with the relative expression profile of the target genes and the FPKM values in the RNA-seq data. The *18s rRNA* and *END04* could be used as reference genes for normalization of the target genes in four cotton species, while *ARF1* was the internal control in YZ1, and *ARF2* and *T2FB* were the target genes in X42 (Fig. 4). A strong, positive correlation coefficient ($R^2=0.8164-0.9984$) was observed between the qRT-PCR results and the RNA RNA-seq data. Moreover, the relative expression profiles and the FPKM values showed similar trends (Fig. 5). Overall, the results indicated the qRT-PCR results were dependable.

Discussion

Analysis of gene expression can successfully predict plant functional genomics²⁵. Because of its high sensitivity, quantitative accuracy and high efficiency, qRT-PCR is also widely used for gene expression studies⁶. However, reference genes may inevitably be influenced by different tissues and treatments, which could lead to unreliable results²⁶. Therefore, it is essential to select valid internal control genes as reference genes for normalization, which can ensure the reliability and accuracy of qRT-PCR data under different experimental conditions⁸. In our study, we conducted high throughout RNA-Seq data on four varieties of *Gossypium hirsutum* L cultivars at three treatment stages (unpublished date) for use as in selecting the reference genes. $|\log_2\text{FoldChange}| < 1$ and a small coefficient of variation (CV) of FPKM of 15 candidate genes were analyzed.

To determine the stability of reference genes expression, we employed three publicly available statistical algorithms, geNorm, NormFinder and BestKeeper. geNorm determined the stability of reference genes by pairwise comparison among test samples¹. In contrast, the NormFinder and BestKeeper algorithms were less sensitive to co-regulation¹⁸. Different studies of the most widely used programs, geNorm and NormFinder, have led to varying conclusions. Many reports showed that they differed slightly, such as in *Oxytropis ochrocephala* Bunge²⁷.

However, other studies indicated that the analytical results obtained using NormFinder were consistent with those obtained with GeNorm, for example, in the halophyte *Halostachys caspica*²⁷, *Chrysanthemum morifolium* and *Chrysanthemum lavandulifolium*²⁸, rice²⁹ and *Rhododendron molle* G. Don³⁰. In this study, the ranking of NormFinder was relatively consistent with analysis of GeNorm. However, the results differed from those obtained using Bestkeeper. For example, for R15, *UBC7* was the most reliable reference gene by obtained by Bestkeeper, while it was ranked at the bottom position in geNorm and NormFinder. The divergence ranking may have been due to the different algorithms³¹. Most studies of Bestkeeper have utilized CV and SD of the Cq to evaluate the stability and expression of reference genes^{30,32}. Besides, as for BestKeeper, which is based on the R of the BestKeeper index by calculating the SD and CV^{30,33,34}. *TAF11* showed the lowest CV± SD values, while the R value was also the lowest in total, YZ1 and X42, to validate results, indicating that *TAF11* would not to be a suitable reference gene. Moreover, *TAF11* was the least reliable gene in geNorm and NormFinder. Thus, the standard of BestKeeper may be an influence factor in the three software programs to calculate stability.

Evaluation using the three specialized software platforms and according to the pairwise variation V value showed that two reference genes were necessary under all conditions. Overall, R15, X33, *ENDO4* and *18s rRNA* were the most stable reference genes. *ENDO4* encodes a putative endonuclease, but no demonstrable endonuclease activity. *18S rRNA* is part of the ribosomal RNA, and thus one of the basic components of all eukaryotic cells. Although ribosomal genes are often suitable reference genes⁹, they are least stable in peaches³⁵. In the present study, *18srRNA* + *ARF1* and *ARF2*+*T2FB* were found to be stable in YZ1 and X42, respectively. *ARF1* and *ARF2* encoded ADP-ribosylation factor, which is essential for vesicle coating and uncoating and functions in GTP-binding^{36,37}. Traditionally, Elongation factor 1- α (*EF1 α*) and Ubiquitin-conjugating enzyme genes are used to as reference genes^{10,38}, while *EF1 α* , *UBC7* and *UBE4* were found to have less stable expression in this study.

To further validate the results of qRT-PCR, the stability of the qRT-PCR was compared with that of RNA-seq data. Using *ARF1* and *18s rRNA* as the internal control for normalization of the target genes in YZ1, *ENDO4* and *18srRNA* as reference genes in R15 and X33, and *ARF2* and *T2FB* as reference genes in X42, which resulted in positive correlation coefficients of $R^2=0.8164-0.9984$. Additionally, the expression pattern of *ERF17* and *AUX22* showed similar trends between the relative expression profile of the target genes and the FPKM values in the RNA-seq data. Similarly, when *18s rRNA* was used as the internal control in the YZ1, *ENDO4* and *18srRNA* normalized the target gene in X42 with a higher R^2 value that was

observed when *ARF2* and *T2FB* were used as reference genes. Additionally, the M value, Sv value, R value and CV \pm SD of the gap between *ENDO4* and *AFR1* was small in YZ1. Similarly, *ENDO4*+*18s rRNA* and *ARF2*+*T2FB* showed a small gap in X42. In summary, *ENDO4* and *18s rRNA* can be used as reliable reference genes in gene expression studies of different conditions and different species in *Gossypium hirsutum*.

Conclusion

We selected suitable reference genes for qRT-PCR normalization in different cotton species among three induction stages. The 15 candidate genes were selected from the transcriptome datasets of *Gossypium hirsutum* L. and assessed by three commonly statistical algorithms, geNorm, NormFinder and BestKeeper. *ENDO4* and *18s rRNA* were identified as appropriate reference genes during three induction stages in four cotton species. The qRT-PCR results and RNA sequencing (RNA-seq) data were strongly positively correlated. The normalized expression profiles of target genes were similar to those of the FPKM values in RNA-seq data for the selected reference genes. Our results indicate that RNA-seq data is a useful source for selecting reference genes in *Gossypium hirsutum*. Additionally, the selected reference genes will provide useful information for appropriate qRT-PCR data normalization in gene expression studies of *Gossypium hirsutum*.

Figure legends

Figure 1. Heat map of candidate reference genes at each pairwise comparison.

Figure 2. Cq values of 15 candidate genes across all of samples.

Figure 3. Average expression stability values (M) and pairwise variation (V_n/V_{n+1}) values of reference genes were calculated by geNorm. a–b: in total, c–d: in YZ1, e–f: in R15, g–h: in X33, i–j: in X42.

Figure 4. qRT-PCR results were validated by comparison with RNA-Seq expression profiles. a: Correlation analysis when *ARF1* and *18s rRNA* were used as the internal controls for normalization of target genes in YZ1, b–c, d–e: *ENDO4* and *18srRNA* as reference genes in R15, X33, YZ1 and X42, respectively, c: *ARF2* and *T2FB* as reference genes in X42.

Figure 5. Validation of the reference gene based on the relative expression level of *AUX22* and *ARF17*. The results are shown as mean fold changes in relative expression when compared to oh. a–c, e–g: The expression profiles of target genes (*AUX22* and *ERF17*) normalized by different stable reference genes as the internal control. d and h: The FPKM values of *AUX22* and *ERF17* in RNA-seq data.

Table 1. Details describing candidate reference genes, primer sequences and amplicon characteristics from qRT-PCR in *Gossypium hirsutum* L.

286 Table 2. The raw value of Cq.
 287 Table 3. Expression stability of the reference genes calculated by GeNorm (M).
 288 Table 4. Expression stability of the reference genes calculated by and NormFinder (Sv).
 289 Table 5. Gene expression stability ranked by Bestkeeper.

290 Acknowledgments

291 This work was supported by National Key R&D Program of China 2017YFD0101604, National Transgenic
 292 Major Projects 2016ZX08005-005, Specific Project for Crops Breeding of Shihezi University (Grant No.
 293 gxjs-yz03)

294 Author contributions

295 A.C., Y.Z., and B.C. conceived and designed the experiments. A.C., D.S. and J.S. performed the experiments.
 296 A.C. and J.S. analyzed the data. A.C., D.S. and Y.Z. contributed reagents/materials/analysis tools. A.C. wrote
 297 the paper. All authors read and approved the final manuscript.
 298

299 Compliance with Ethical Standards

300 The authors declare that they have no conflicts of interest associated with this paper.

301 References

- 302 1 Vandesompele, J. *et al.* Accurate normalization of real-time quantitative RT-PCR data by geometric
 303 averaging of multiple internal control genes. *Genome Biology* **3** (2002).
- 304 2 Shi, X. *et al.* De novo comparative transcriptome analysis provides new insights into sucrose induced
 305 somatic embryogenesis in camphor tree (*Cinnamomum camphora* L.). *BMC Genomics* **17**, 26 (2016).
- 306 3 Yang, M. *et al.* Transcriptomic Analysis of the Regulation of Rhizome Formation in Temperate and
 307 Tropical Lotus (*Nelumbo nucifera*). *Scientific Reports* **5**, 13059 (2015).
- 308 4 Bustin, S. A. Quantification of mRNA using real-time reverse transcription PCR (RT-PCR): trends and
 309 problems. *Journal of Molecular Endocrinology* **29**, 23 (2002).
- 310 5 Ma, R., Xu, S., Zhao, Y., Xia, B. & Wang, R. Selection and Validation of Appropriate Reference Genes for
 311 Quantitative Real-Time PCR Analysis of Gene Expression in *Lycoris aurea*. *Front Plant Sci* **7** (2016).
- 312 6 Hong, S. Y., Seo, P. J., Yang, M. S., Xiang, F. & Park, C. M. Exploring valid reference genes for gene
 313 expression studies in *Brachypodium distachyon* by real-time PCR. *BMC Plant Biology* **8**, 112 (2008).
- 314 7 Derveaux, S., Vandesompele, J. & Hellemans, J. How to do successful gene expression analysis using
 315 real-time PCR. *Methods* **50**, 227-230 (2010).
- 316 8 Huggett, J., Dheda, K., Bustin, S. & Zumla, A. Real-time RT-PCR normalisation; strategies and
 317 considerations. *Genes & Immunity* **6**, 279 (2005).

- 318 9 Zhang, S., Zeng, Y., Yi, X. & Zhang, Y. Selection of suitable reference genes for quantitative RT-PCR
319 normalization in the halophyte *Halostachys caspica* under salt and drought stress. *Scientific Reports* **6**,
320 30363 (2016).
- 321 10 Zhu, J. *et al.* Reference Gene Selection for Quantitative Real-time PCR Normalization in *Caragana*
322 *intermedia* under Different Abiotic Stress Conditions. *Plos One* **8**, e53196 (2013).
- 323 11 Sun, H. P., Li, F., Ruan, Q. M. & Zhong, X. H. Identification and validation of reference genes for
324 quantitative real-time PCR studies in *Hedera helix* L. *Plant Physiology & Biochemistry* **108**, 286-294
325 (2016).
- 326 12 Kumria, R. *et al.* High-frequency somatic embryo production and maturation into normal plants in
327 cotton (*Gossypium hirsutum*) through metabolic stress. *Plant Cell Reports* **21**, 635-639 (2003).
- 328 13 Yang, X. *et al.* Transcript profiling reveals complex auxin signalling pathway and transcription
329 regulation involved in dedifferentiation and redifferentiation during somatic embryogenesis in
330 cotton. *BMC Plant Biology* **12**, 110 (2012).
- 331 14 Cheng, W. H. *et al.* De novo transcriptome analysis reveals insights into dynamic homeostasis
332 regulation of somatic embryogenesis in upland cotton (*G. hirsutum* L.). *Plant Molecular Biology* **92**,
333 279-292 (2016).
- 334 15 Yang, X. Y. & Zhang, X. L. Regulation of somatic embryogenesis in higher plants. *Critical Reviews in*
335 *Plant Sciences* **29**, 36-57 (2010).
- 336 16 Radonić, A. *et al.* Guideline to reference gene selection for quantitative real-time PCR. *Biochemical &*
337 *Biophysical Research Communications* **313**, 856 (2004).
- 338 17 Andersen, C. L., Jensen, J. L. & TF, Ø. Normalization of real-time quantitative reverse transcription-
339 PCR data: a model-based variance estimation approach to identify genes suited for normalization,
340 applied to bladder and colon cancer data sets. *Cancer research* **64**, 5245 (2004).
- 341 18 Pfaffl, M. W., Tichopad, A., Prgomet, C. & Neuvians, T. P. Determination of stable housekeeping genes,
342 differentially regulated target genes and sample integrity: BestKeeper – Excel-based tool using pair-
343 wise correlations. *Biotechnology Letters* **26**, 509-515 (2004).
- 344 19 Jin, S., Zhang, X., Y, Guo, X., Liang, S. & Zhu, H. Identification of a novel elite genotype for in vitro
345 culture and genetic transformation of cotton. *Biologia Plantarum* **50**, 519-524 (2006).
- 346 20 Cao, A. *et al.* Comparative Transcriptome Analysis of SE initial dedifferentiation in cotton of different
347 SE capability. *Scientific Reports* **7** (2017).
- 348 21 Jonge, H. J. M. D. *et al.* Evidence Based Selection of Housekeeping Genes. *Plos One* **2**, : e898. (2007).
- 349 22 Bustin, S. A. Why the need for qPCR publication guidelines?—The case for MIQE. *Methods* **50**, 217
350 (2010).
- 351 23 Livak, K. J. & Schmittgen, T. D. Analysis of relative gene expression data using real-time quantitative
352 PCR and the 2(-Delta Delta C(T)) Method. *Methods* **25**, 402-408 (2001).
- 353 24 Chen, Y. *et al.* Selection of reference genes for quantitative real-time PCR normalization in creeping
354 bentgrass involved in four abiotic stresses. *Plant Cell Reports* **34**, 1825 (2015).
- 355 25 Li, C. *et al.* An improved fruit transcriptome and the identification of the candidate genes involved in
356 fruit abscission induced by carbohydrate stress in litchi. *Frontiers in Plant Science* **6**, 439 (2015).
- 357 26 Dheda, K. *et al.* The implications of using an inappropriate reference gene for real-time reverse
358 transcription PCR data normalization. *Analytical biochemistry* **344**, 141 (2005).

- 27 Zhuang, H., Fu, Y., He, W., Wang, L. & Wei, Y. Selection of appropriate reference genes for quantitative real-time PCR in *Oxytropis ochrocephala* Bunge using transcriptome datasets under abiotic stress treatments. *Frontiers in Plant Science* **6**, 475 (2015).
- 28 Qi, S. *et al.* Reference Gene Selection for RT-qPCR Analysis of Flower Development in *Chrysanthemum morifolium* and *Chrysanthemum lavandulifolium*. *Frontiers in Plant Science* **7** (2016).
- 29 Jain, M. Genome-wide identification of novel internal control genes for normalization of gene expression during various stages of development in rice. *Plant Science* **176**, 702-706 (2009).
- 30 Xiao, Z. *et al.* Selection of Reliable Reference Genes for Gene Expression Studies on *Rhododendron molle* G. Don. *Frontiers in Plant Science* **7**, 1547 (2016).
- 31 Niu, X. *et al.* Selection of reliable reference genes for quantitative real-time PCR gene expression analysis in Jute (*Corchorus capsularis*) under stress treatments. *Frontiers in Plant Science* **6**, 848 (2015).
- 32 Zhao, Y. *et al.* Selection of Reference Genes for Gene Expression Normalization in *Peucedanum praeruptorum* Dunn under Abiotic Stresses, Hormone Treatments and Different Tissues. *Plos One* **11**, e0152356 (2016).
- 33 Kong, F., Cao, M., Sun, P., Liu, W. & Mao, Y. Selection of reference genes for gene expression normalization in *Pyropia yezoensis* using quantitative real-time PCR. *Journal of Applied Phycology* **27**, 1003-1010 (2015).
- 34 Shi, X. *et al.* De novo comparative transcriptome analysis provides new insights into sucrose induced somatic embryogenesis in camphor tree (*Cinnamomum camphora* L.). *Bmc Genomics* **17**, 26 (2016).
- 35 Tong, Z., Gao, Z., Wang, F., Zhou, J. & Zhang, Z. Selection of reliable reference genes for gene expression studies in peach using real-time PCR. *BMC Molecular Biology* **10**, 71 (2009).
- 36 Matheson, L. A. *et al.* Multiple Roles of ADP-Ribosylation Factor 1 in Plant Cells Include Spatially Regulated Recruitment of Coatomer and Elements of the Golgi Matrix. *Plant Physiology* **143**, 1615-1627 (2007).
- 37 Serventi, I. M., Cavanaugh, E., Moss, J. & Vaughan, M. Characterization of the gene for ADP-ribosylation factor (ARF) 2, a developmentally regulated, selectively expressed member of the ARF family of approximately 20-kDa guanine nucleotide-binding proteins. *The Journal of biological chemistry* **268**, 4863-4872 (1993).
- 38 Chang, E. *et al.* Selection of reference genes for quantitative gene expression studies in *Platycladus orientalis* (Cupressaceae) Using real-time PCR. *Plos One* **7**, e33278 (2011).

Figure 1

Figure 1

Heat map of candidate reference genes at each pairwise comparison

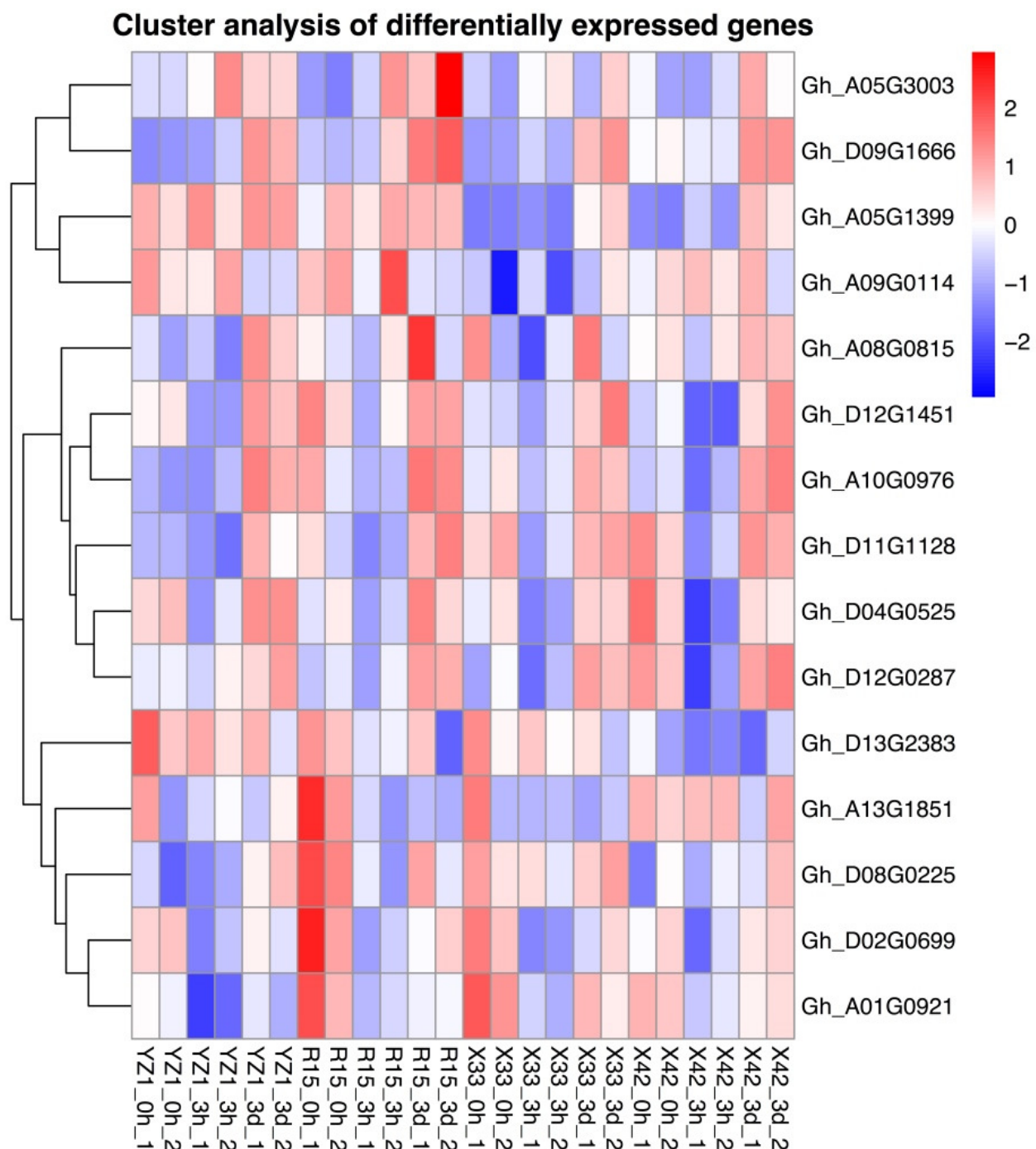


Figure 2(on next page)

figure2

Cq values of 15 candidate genes across all of samples

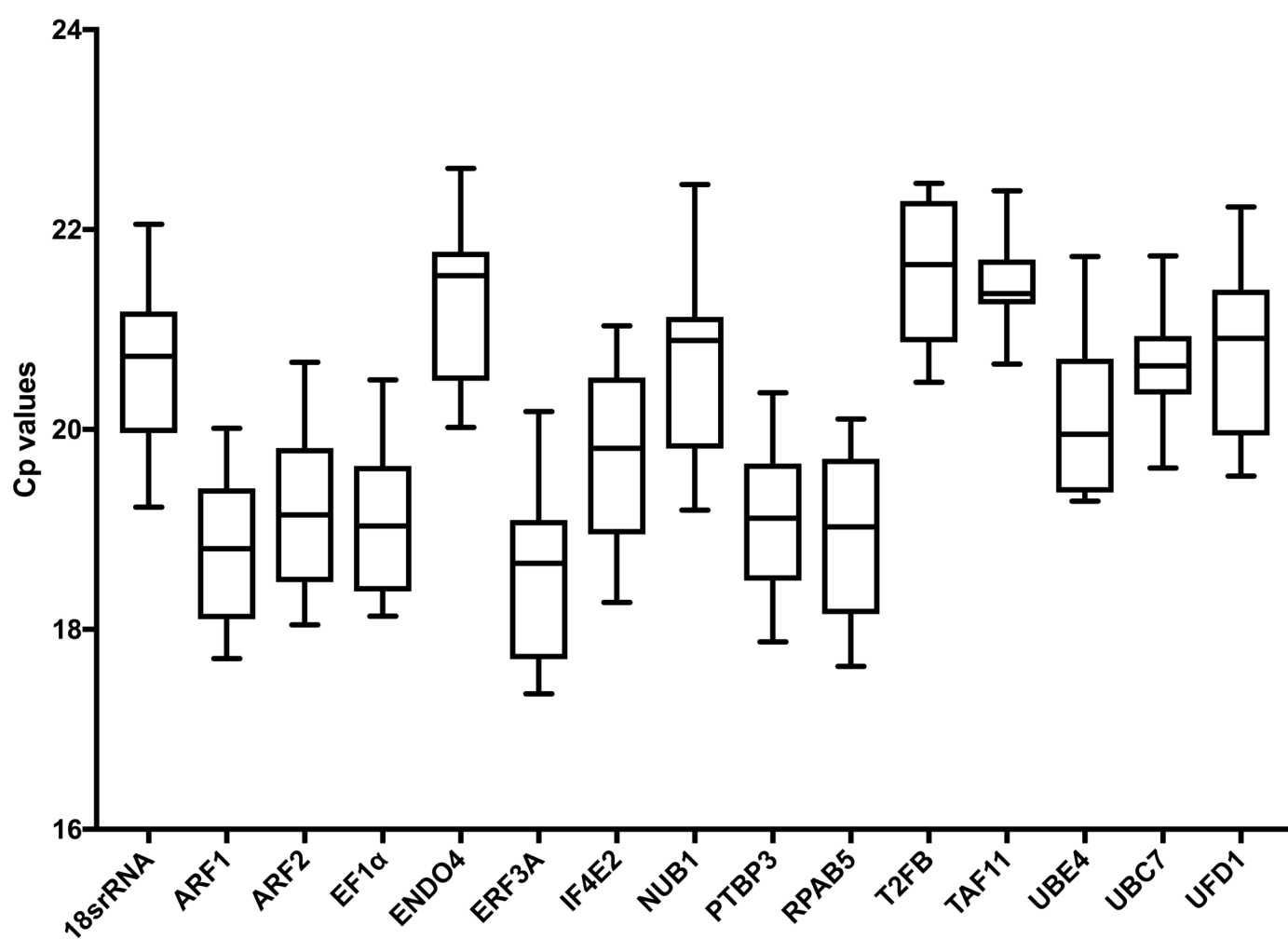


Figure 3

Figure 3

Average expression stability values (M) and pairwise variation (V_n/V_{n+1}) values of reference genes were calculated by geNorm. a-b: in total, c-d: in YZ1, e-f: in R15, g-h: in X33, i-j: in X42.

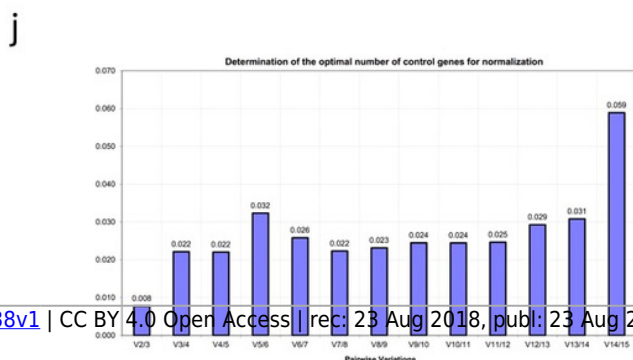
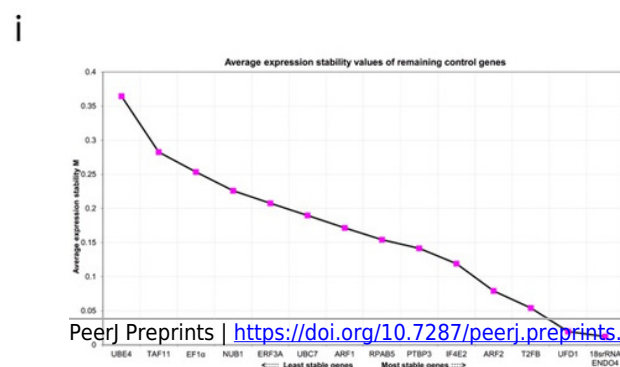
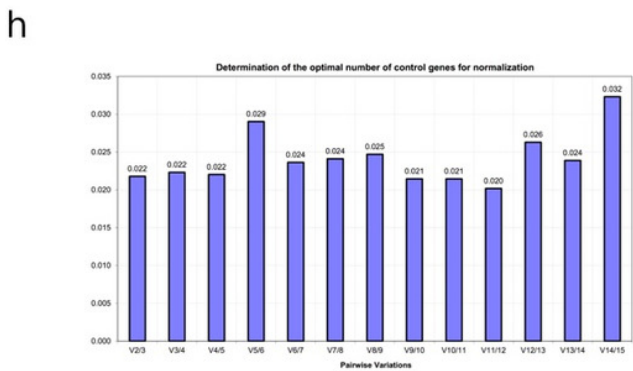
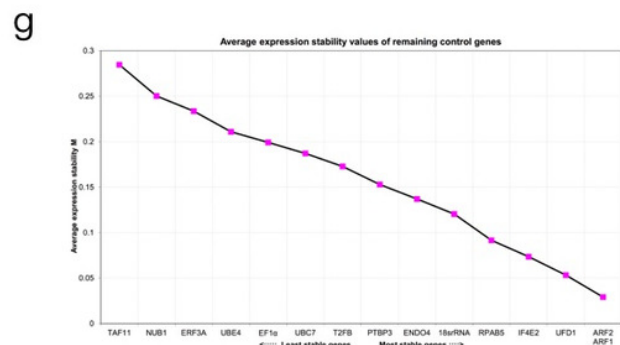
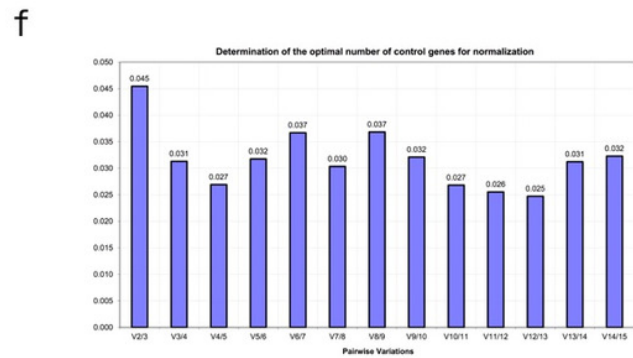
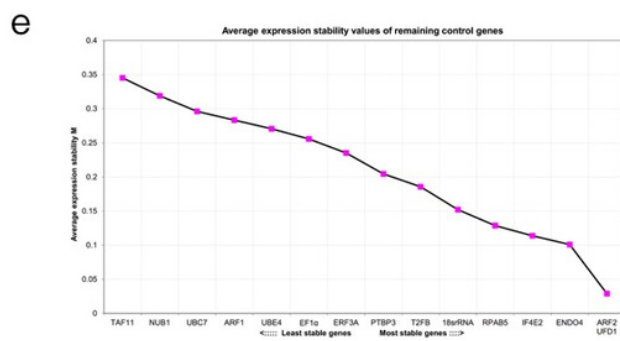
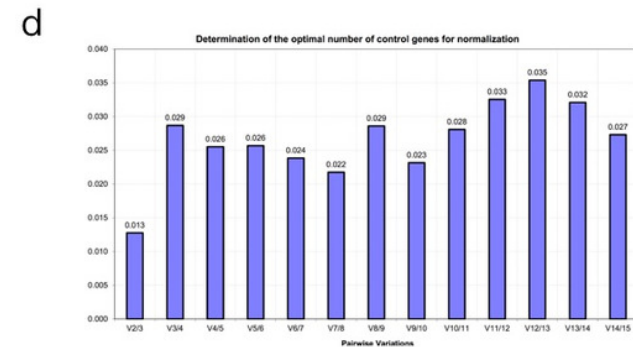
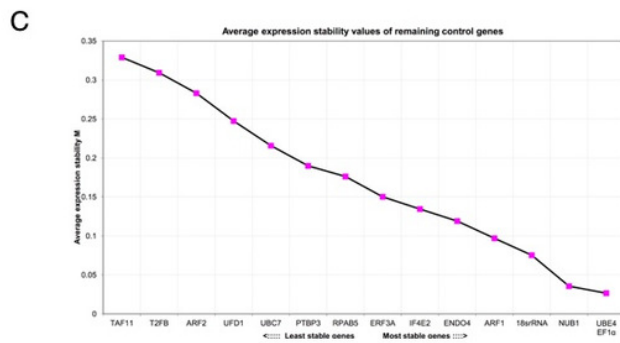
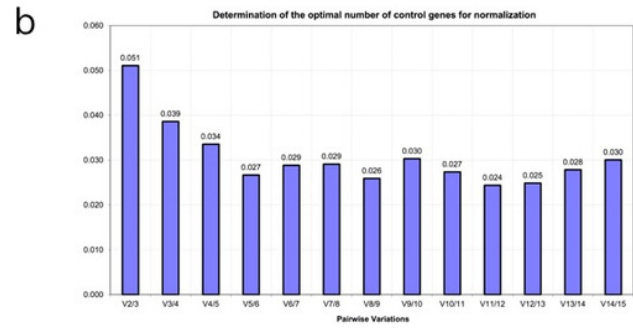
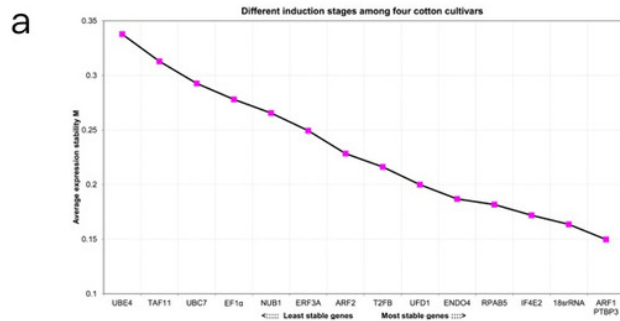


Figure 4 (on next page)

Figure 4

qRT-PCR results were validated by comparison with RNA-Seq expression profiles. a: Correlation analysis when ARF1 and 18s rRNA were used as the internal controls for normalization of target genes in YZ1, b-c, d-e: ENDO4 and 18srRNA as reference genes in R15, X33 , YZ1 and X42, respectively, c: ARF2 and T2FB as reference genes in X42.

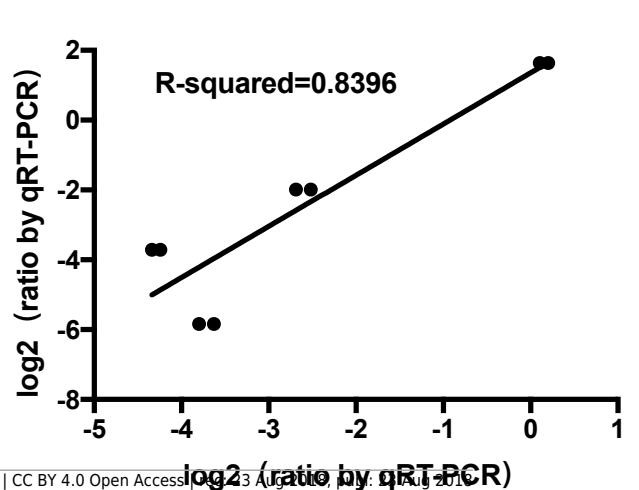
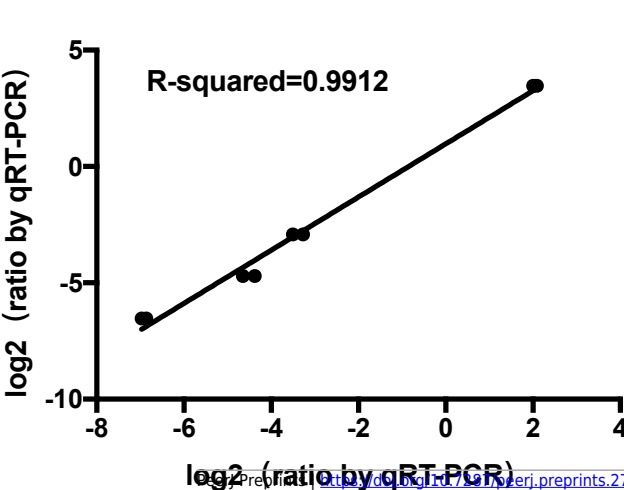
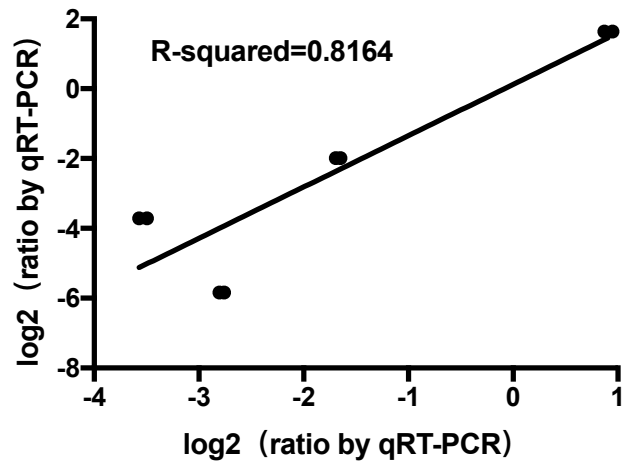
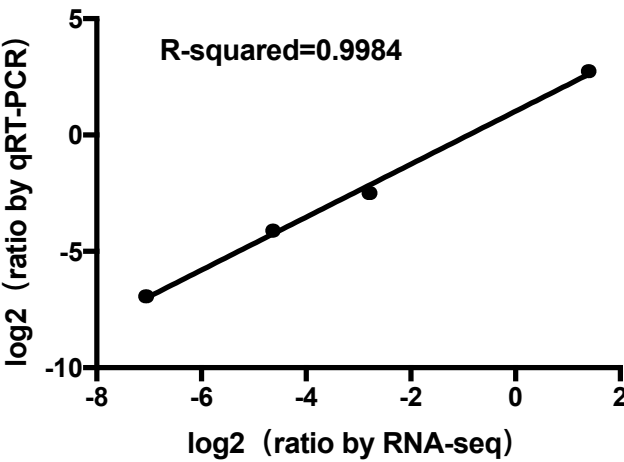
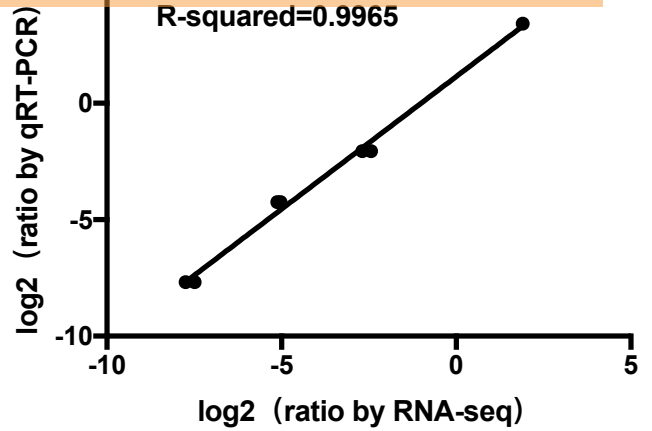
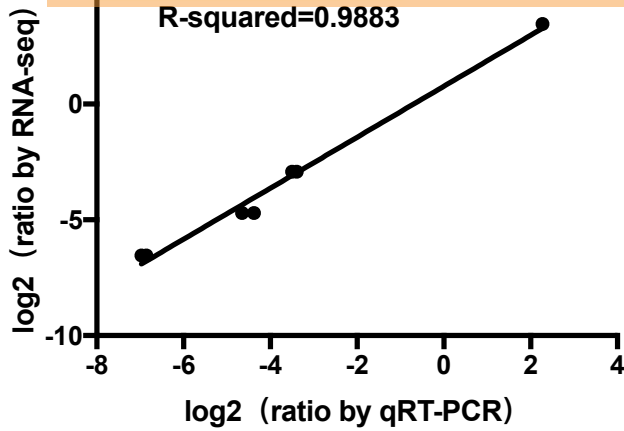


Figure 5(on next page)

Figure 5

Validation of the reference gene based on the relative expression level of AUX22 and ARF17. The results are shown as mean fold changes in relative expression when compared to oh. a-c, e-g: The expression profiles of target genes (AUX22 and ERF17) normalized by different stable reference genes as the internal control. d and h: The FPKM values of AUX22 and ERF17 in RNA-seq data.

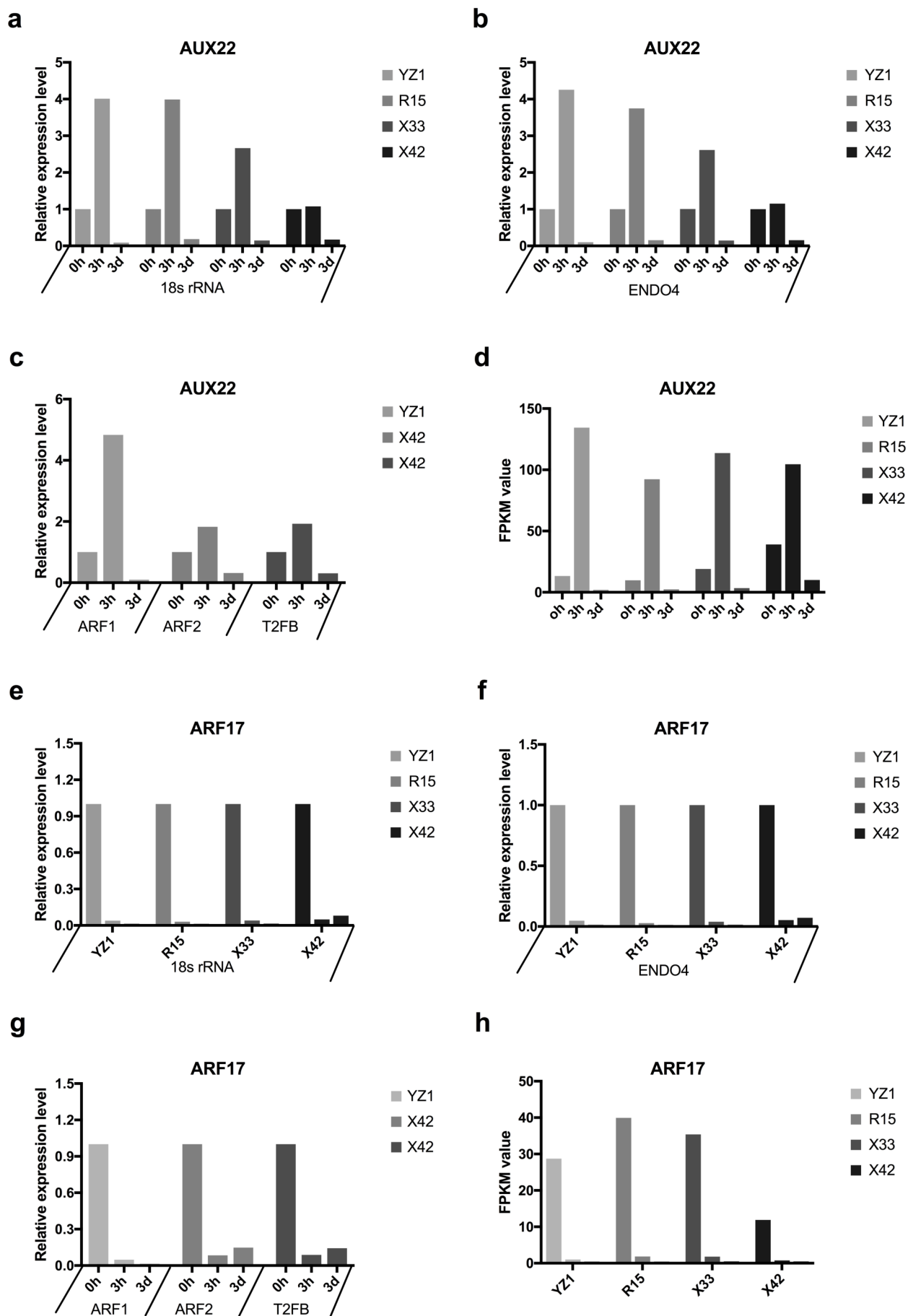


Table 1 (on next page)

Table 1

Details describing candidate reference genes, primer sequences and amplicon characteristics from qRT-PCR in *Gossypium hirsutum* L.

Table 1. Details of candidate reference genes, primer sequences and amplicon characteristics from qRT-PCR in *Gossypium hirsutum* L.

Gene name	Primer sequence (5'-3')	Length (bp)	E (%)	R2
18srRNA	TTACGCAATGCGCTCTGGA ACCGCAGAGCTGACAGATG	117	104.70	0.9968
ARF1	CTGTGAGCAGAAAGTGGAAGC CAGCTGCATCAAGACCCACC	111	104.61	0.9996
ARF2	CCACTTCTGGTGAAGGTCTGT AACTCTAAAAGGGGCCAGCA	118	100.33	0.9980
EF1 α	CAGCTTCAGATCGCTTCTATTTCT TGGCCAGTGGTGGTTGACTT	124	100.07	0.9998
ENDO4	TTGACAGAGGCGCTGATGTT CTGCGGTACCAACTGACTGT	116	100.44	0.9905
ERF3A	GCCCTATTTGCCACAAAACCC TTCAGGAATGAGCGTGGCAT	140	101.56	0.9886
IF4E2	CAAGACTGCAACGAATGAGGC GCTCAAACATTGTATCGACCTTTCA	144	100.36	0.9894
NUB1	TTGCACTACATATGAGGTTGGAGTT AGGCTTCATCAGGCACTTGTA	132	101.60	0.9677
PTBP3	GTCCTTGCAAATGGCGGAAG CCTGATTCTTTGCACGGAGC	140	102.08	0.9961
RPAB5	CTTCACCTGCGGAAAGGTCA AGCAGTACCGAACCAATCCC	113	99.85	1.0643
T2FB	GGATCGCGGGGAATTGGAA TGCCTCTCTTATTGTACACGCA	149	99.09	1.0030
TAF11	TCGTCTGCATTAGAGAGTCGC GGCTGTTTCTAGCTCATCCTCA	138	96.13	0.9999
UBC7	GCTGGACGCCAGTACATACA GGCTGACCTTCCTCTGAAT	144	99.21	1.0997
UBE4	TGGGCCCCCTTTTCCATGTT TCAGCTGCTCGTCTAGTTGATG	109	120.30	0.9999
UFD1	TGTCAGCCGTTCTAAGGAAACA ACTTTCTCCCGGTGAATGGC	107	115.58	0.9773

4
5
6

Table 2 (on next page)

Table 2

The raw value of Cq.

1

Table 2. Cq value of candidate reference genes

	YZ1- 0h	YZ1- 3h	YZ1- 3d	R15- 0h	R15- 3h	R15- 3d	X33- 0h	X33- 3h	X33- 3d	X42- 0h	X42- 3h	X42- 3d
18srRNA	20.98	20.42	19.99	21.72	20.72	19.95	22.05	21.24	19.89	20.75	21.02	19.22
ARF1	19.01	18.48	18.02	19.76	18.90	18.35	20.01	19.47	17.96	18.72	19.24	17.71
ARF2	18.70	18.90	18.58	20.67	19.50	18.37	20.41	19.89	18.44	19.39	19.60	18.05
EF1 α	19.68	18.80	18.47	20.50	18.85	18.13	20.23	19.34	18.35	19.51	19.22	18.16
ENDO4	21.56	21.10	20.41	22.61	21.52	20.59	22.38	21.65	20.45	21.59	21.82	20.02
ERF3A	19.15	18.34	17.70	20.04	18.47	17.63	20.18	18.93	17.72	18.95	18.85	17.36
IF4E2	20.05	19.65	18.94	20.93	19.94	19.00	21.04	20.59	18.89	19.69	20.32	18.27
NUB1	21.05	20.21	19.89	22.45	20.76	19.72	22.29	21.02	19.79	21.14	21.10	19.19
PTBP3	19.19	19.04	18.40	20.01	19.32	18.49	20.37	19.60	18.49	19.02	19.68	17.88
RPAB5	19.05	18.84	18.32	20.11	19.25	18.10	19.98	19.69	18.01	19.01	19.71	17.63
T2FB	21.29	21.57	21.13	22.46	21.97	20.79	22.46	22.33	20.76	21.73	22.16	20.47
TAF11	21.47	21.15	21.33	22.39	21.29	21.31	22.38	21.74	21.24	21.38	21.57	20.66
UBE4	20.64	19.83	19.42	21.66	20.07	19.64	21.73	20.73	19.35	20.39	19.32	19.28
UBC7	20.82	20.43	20.35	21.74	20.74	20.35	21.23	20.82	19.66	20.53	20.98	19.61
UFD1	20.57	20.82	19.94	22.23	21.01	19.94	22.01	21.42	19.85	21.02	21.32	19.54

2

3

4

5

6

Table 3(on next page)

Table 3

Expression stability of the reference genes calculated by GeNorm (M).

1

Table 3. Expression stability of the reference genes calculated by GeNorm (M)

Raw	Gene	all	Gene	YZ1	Gene	R15	Gene	X33	Gene	X42
1	18srRNA	0.27	ARF1	0.24	ENDO4	0.25	ARF2	0.21	ARF2	0.26
2	ENDO4	0.27	18srRNA	0.25	IF4E2	0.26	ARF1	0.21	18srRNA	0.26
3	ARF1	0.29	RPAB5	0.26	18srRNA	0.27	ENDO4	0.21	UFD1	0.26
4	PTBP3	0.29	ENDO4	0.26	UFD1	0.29	18srRNA	0.22	T2FB	0.26
5	RPAB5	0.30	PTBP3	0.27	RPAB5	0.30	UFD1	0.23	ENDO4	0.27
6	IF4E2	0.30	IF4E2	0.28	ARF2	0.31	PTBP3	0.23	ARF1	0.31
7	UFD1	0.30	NUB1	0.29	PTBP3	0.34	UBE4	0.27	PTBP3	0.32
8	ARF2	0.34	UBE4	0.31	UBE4	0.35	EF1 α	0.27	ERF3A	0.32
9	ERF3A	0.34	EF1 α	0.32	ARF1	0.35	IF4E2	0.27	UBC7	0.33
10	T2FB	0.35	UBC7	0.33	ERF3A	0.35	RPAB5	0.28	IF4E2	0.34
11	EF1 α	0.35	ERF3A	0.38	EF1 α	0.36	UBC7	0.30	RPAB5	0.36
12	NUB1	0.38	UFD1	0.40	UBC7	0.37	T2FB	0.32	NUB1	0.40
13	UBC7	0.38	ARF2	0.43	T2FB	0.38	ERF3A	0.36	EF1 α	0.41
14	TAF11	0.45	T2FB	0.45	NUB1	0.48	NUB1	0.37	TAF11	0.47
15	UBE4	0.50	TAF11	0.46	TAF11	0.52	TAF11	0.51	UBE4	0.90

2

3

4

5

Table 4(on next page)

Table 4

Expression stability of the reference genes calculated by and NormFinder (Sv).

1

Table 4. Expression stability of the reference genes calculated by NormFinder (Sv)

Rank	Gene	All	Gene	YZ1	Gene	R15	Gene	X33	Gene	X42
1	18srRNA	0.08	RPAB5	0.04	ENDO4	0.01	ENDO4	0.02	ARF2	0.04
2	ENDO4	0.09	ARF1	0.05	IF4E2	0.03	ARF2	0.05	T2FB	0.06
3	ARF1	0.10	18srRNA	0.06	18srRNA	0.05	18srRNA	0.05	UFD1	0.08
4	PTBP3	0.11	PTBP3	0.08	UFD1	0.13	PTBP3	0.06	18srRNA	0.09
5	IF4E2	0.13	IF4E2	0.11	RPAB5	0.13	ARF1	0.06	ENDO4	0.10
6	RPAB5	0.13	ENDO4	0.12	ARF2	0.14	UFD1	0.09	ERF3A	0.11
7	UFD1	0.13	UBC7	0.15	UBE4	0.17	EF1 α	0.11	ARF1	0.12
8	ARF2	0.17	UBE4	0.17	PTBP3	0.18	IF4E2	0.14	UBC7	0.14
9	ERF3A	0.17	NUB1	0.17	ARF1	0.19	UBE4	0.15	PTBP3	0.16
10	EF1 α	0.18	EF1 α	0.19	ERF3A	0.19	UBC7	0.15	EF1 α	0.20
11	T2FB	0.18	UFD1	0.23	EF1 α	0.20	RPAB5	0.16	IF4E2	0.20
12	UBC7	0.21	ERF3A	0.23	UBC7	0.20	T2FB	0.20	RPAB5	0.22
13	NUB1	0.21	ARF2	0.27	T2FB	0.22	ERF3A	0.22	NUB1	0.23
14	TAF11	0.27	T2FB	0.28	NUB1	0.32	NUB1	0.23	TAF11	0.25
15	UBE4	0.31	TAF11	0.28	TAF11	0.35	TAF11	0.35	UBE4	0.60

2

Table 5(on next page)

Table 5

Gene expression stability ranked by Bestkeeper.

1

Table 5. Gene expression stability ranked by Bestkeeper

	Total	YZ1	R15	X33	X42
Gene (R)	TAF11 (0.94)	TAF11 (0.46)	TAF11 (0.95)	TAF11 (0.98)	TAF11 (0.52)
CV±SD	1.63±0.35	0.53±0.11	2.23±0.48	1.81±0.39	1.72±0.36
Gene (R)	UBC7 (0.94)	ARF2 (0.44)	UBC7 (1.00)	UBC7 (0.99)	UBE4 (0.98)
CV±SD	2.18±0.45	0.62±0.12	2.52±0.53	2.94±0.61	2.46±0.48
Gene (R)	T2FB (0.95)	T2FB (0.44)	ARF1 (1.00)	END04 (1.00)	UBC7 (0.95)
CV±SD	2.74±0.59	0.75±0.16	2.65±0.50	3.23±0.69	2.49±0.51
Gene (R)	PTBP3 (0.98)	UBC7 (0.94)	PTBP3 (1.00)	EF1α (0.99)	EF1α (0.93)
CV±SD	2.99±0.57	0.94±0.19	2.72±0.52	3.29±0.64	2.83±0.54
Gene (R)	18srRNA(0.98)	RPAB5 (0.98)	T2FB (0.96)	T2FB (0.97)	ARF1 (0.99)
CV±SD	3.09±0.64	1.48±0.28	2.91±0.63	3.33±0.73	3.05±0.57
Gene (R)	ENDO4 (0.99)	UFD1 (0.67)	18srRNA (0.99)	PTBP3 (1.00)	T2FB (1.00)
CV±SD	3.11±0.66	1.63±0.33	2.97±0.62	3.40±0.66	3.05±0.65
Gene (R)	ARF1 (0.98)	PTBP3 (0.97)	END04 (1.00)	18srRNA (1.00)	ARF2 (1.00)
CV±SD	3.17±0.60	1.67±0.32	3.21±0.69	3.70±0.78	3.39±0.64
Gene (R)	UFD1 (0.96)	18srRNA(0.99)	IF4E2 (1.00)	ARF2 (1.00)	PTBP3 (0.98)
CV±SD	3.35±0.70	1.69±0.35	3.25±0.65	3.89±0.76	3.47±0.65
Gene (R)	EF1α (0.95)	ARF1 (0.99)	RPAB5 (0.99)	UFD1 (1.00)	UFD1 (0.99)
CV±SD	3.36±0.64	1.82±0.34	3.66±0.70	3.92±0.83	3.53±0.73
Gene (R)	RPAB5 (0.98)	END04 (1.00)	UFD1 (0.67)	NUB1 (0.98)	ENDO4 (1.00)
CV±SD	3.49±0.66	1.95±0.41	3.70±0.78	3.98±0.84	3.54±0.75
Gene (R)	UBE4 (0.92)	IF4E2 (1.00)	UBE4 (0.99)	UBE4 (1.00)	18srRNA(0.99)
CV±SD	3.54±0.71	2.08±0.41	3.91±0.80	4.05±0.83	3.63±0.74
Gene (R)	IF4E2 (0.99)	NUB1 (0.95)	ARF2 (1.00)	ARF1 (1.00)	ERF3A (0.98)
CV±SD	3.55±0.70	2.18±0.44	3.96±0.77	4.12±0.79	3.73±0.69
Gene (R)	ARF2 (0.96)	UBE4 (0.99)	EF1α (0.99)	RPAB5 (0.98)	IF4E2 (0.99)
CV±SD	3.66±0.70	2.24±0.45	4.66±0.89	4.22±0.81	3.97±0.77
Gene (R)	NUB1 (0.98)	EF1α (0.94)	NUB1 (1.00)	IF4E2 (0.99)	RPAB5 (0.99)
CV±SD	3.85±0.80	2.44±0.46	4.69±0.98	4.24±0.85	4.09±0.77
Gene (R)	ERF3A (0.98)	ERF3A (0.99)	ERF3A (0.99)	ERF3A (0.98)	NUB1 (0.98)
CV±SD	3.98±0.74	2.72±0.50	4.74±0.89	4.35±0.82	4.18±0.86

2

SD: standard deviation of the Cq; CV: coefficient of variance expression as a percentage of the Cq level; R:

3

correlation coefficient of the Cq level.

