# Genome-wide identification and analysis of cystatin family genes in Sorghum (*Sorghum bicolor* L.)

Jie Li [Corresp., 1] , Xinhao Liu [2] , Jingmei Wang [2] , Junyan Sun [1] , Dexian He [Corresp. 3]

[1] College of Agronomy, Xinyang Agriculture and Forestry University, Xinyang, Henan Province, China

[2] Central Laboratory, Xinyang Agriculture and Forestry University, Xinyang, Henan Province, China

[3] Collaborative Innovation Center of Henan Grain Crops/National Key Laboratory of Wheat and Maize Crop Science, College of Agronomy, Henan Agricultural University, Zhengzhou, China

Corresponding Authors: Jie Li, Dexian He
Email address: ljljlxh123@126.com, hedexian@126.com

To set a systematic study of the Sorghum *cystatins* (*SbCys*) gene family, a comprehensive genome-wide analysis of the *SbCys* family genes was performed by bioinformatics-based methods. In total, 18 *SbCys* genes were identified in Sorghum, which distributed unevenly on chromosomes, and two genes were involved in tandem duplication event. All *SbCys* genes had similar exon/intron structure and motifs, indicating their high evolutionary conservation. Transcriptome analysis showed that 16 *SbCys* genes were expressed in at least one tested tissues, and most genes displayed higher expression levels in reproductive tissues than in vegetable tissues, indicating that the *SbCys* genes participated in the regulation of seed formation. Furthermore, the expressions of 7 *SbCys* genes were induced by *Bipolaris sorghicola* infection, while only 2 genes were responsive to aphid infestation. In addition, quantitative real-time polymerase chain reaction (qRT-PCR) confirmed that 17 *SbCys* genes were induced by one or two abiotic stresses (dehydration, salt shock and ABA stresses). In addition, the interaction network indicated that SbCys proteins were associated with several biological processes, including seed development and stress responses. Notably, the expression of *SbCys4* was up-regulated under biotic and abiotic stresses, suggesting its potential roles in mediating the response of Sorghum to adverse environmental impact. Our results provide new insights into the structural and functional characteristics of *SbCys* gene family, which laid the foundation for better understanding the roles and regulatory mechanism of Sorghum cystatins in seed development and responses to different stress conditions.

1 # Genome-wide identification and analysis of cystatin family genes

2 # in Sorghum (*Sorghum bicolor* L.)

3 Jie Li[1], Xinhao Liu[2], Jingmei Wang[2], Junyan Sun[1], Dexian He[3*]

4

5 **Author affiliation:**

6 1 College of Agronomy, Xinyang Agriculture and Forestry University, Xinyang, Henan 464001,

7 China

8 2 Central Laboratory, Xinyang Agriculture and Forestry University, Xinyang, Henan 464001,

9 China

10 3 Collaborative Innovation Center of Henan Grain Crops/National Key Laboratory of Wheat and

11 Maize Crop Science, College of Agronomy, Henan Agricultural University, Zhengzhou, Henan,

12 450002, China

13

14 [*]Corresponding author: Dexian He

15 College of Agronomy, Henan Agricultural University, Zhengzhou, Henan, 450002, China

16 E-mail: hedexian@126.com

17

18

19

20

21

22

23

24  **ABSTRACT**

25  To set a systematic study of the Sorghum *cystatins* (*SbCys*) gene family, a comprehensive

26  genome-wide analysis of the *SbCys* family genes was performed by bioinformatics-based

27  methods. In total, 18 *SbCys* genes were identified in Sorghum, which distributed unevenly on

28  chromosomes, and two genes were involved in tandem duplication event. All *SbCys* genes had

29  similar exon/intron structure and motifs, indicating their high evolutionary conservation.

30  Transcriptome analysis showed that 16 *SbCys* genes were expressed in at least one tested tissues,

31  and most genes displayed higher expression levels in reproductive tissues than in vegetable

32  tissues, indicating that the *SbCys* genes participated in the regulation of seed formation.

33  Furthermore, the expressions of 7 *SbCys* genes were induced by *Bipolaris sorghicola* infection,

34  while only 2 genes were responsive to aphid infestation. In addition, quantitative real-time

35  polymerase chain reaction (qRT-PCR) confirmed that 17 *SbCys* genes were induced by one or

36  two abiotic stresses (dehydration, salt shock and ABA stresses). In addition, the interaction

37  network indicated that SbCys proteins were associated with several biological processes,

38  including seed development and stress responses. Notably, the expression of *SbCys4* was up-

39  regulated under biotic and abiotic stresses, suggesting its potential roles in mediating the

40  response of Sorghum to adverse environmental impact. Our results provide new insights into the

41  structural and functional characteristics of *SbCys* gene family, which laid the foundation for

42  better understanding the roles and regulatory mechanism of Sorghum cystatins in seed

43  development and responses to different stress conditions.

44

45

PeerJ

46

**INTRODUCTION**

47

48  Cystatins are competitive and reversible inhibitors of cysteins proteases from families C1A and

49  C13, which have been identified in many plant species (Martinez and Diaz, 2008; Zhao et al.

50  2014). Based on their primary sequence homology, three signature motifs include a QxVxG

51  reactive site, a tryptophan residue (W) located downstream of the reactive site, and one or two

52  glycine (G) residues in the flexible N terminus of the protein. These three motifs are important

53  for the cystatin inhibitory mechanism (Jenko et al. 2003; Stubbs et al. 1990). In addition, a

54  consensus sequence ([LVI]-[AGT]-[RKE]-[FY]-[AS]-[VI]-x-[EDQV]-[HYFQ]-N) in cystatins is

55  conformed to a predicted secondary α-helix structure (Margis et al. 1998). Most plant cystatins

56  are small proteins with a molecular mass in the 12- to 16-kD range (Margis et al. 1998). Some

57  plant cystatins contain a C-terminal extension that raises their molecular weights up to 23 kDa,

58  are thought to be involved in the inhibition of cysteine protease activities in the peptidase C13

59  family (Martinez et al. 2007; Martinez and Diaz, 2008).

60  The principal functions of plant cystatins are related to the regulation of endogenous

61  cysproteases during plant growth and development, senescence, and programmed cell death

62  (Belenghi et al. 2010; Díazmendoza et al. 2014; Zhao et al. 2014). Additionally, Plant cystatins

63  have been used as effective molecules against different pests and pathogens (Martinez et al.

64  2016). For example, several publications reported the inhibition of recombinant cystatins on the

65  growth of some pests and fungi (Martinez et al. 2005; Lima et al. 2015). Tomato plants over-

66  expressing the wheat cystatin *TaMDC1* displayed a broad stress resistance for bacterial pathogen,

67  and the defense responses were mediated by methyl jasmonate and salicylic acid (Christova et al.

68  2018). The inhibition of amaranth cystatin on the digestive insect cysteine endopeptidases was

69    observed by Valdés-Rodríguez et al. (2015). Plant cystatins are also involved in the responses to

70    abiotic stresses, such as over-expression of *MpCYS4* in apple delayed natural and stress-induced

71    leaf senescence (Tan et al. 2017). Song et al. (2017) found that the expression of *AtCYS5* was

72    induced by heat stress (HS) and exogenous ABA treatment in germinating seed, furthermore,

73    over expression of *AtCYS5* enhanced HS tolerance in transgenic *Arabidopsis*.

74    To date, plant cystatin family genes had been well described in several plant species such as

75    *Arabidopsis*, rice, soybean, wheat, and *Populus trichocarpa* (Martinez and Diaz, 2008; Wang et

76    al. 2015; Yuan et al. 2016; Dutt et al. 2016). However, a genome-wide study of cystatins family

77    genes in Sorghum (*Sorghum bicolor* L.) has not yet been performed. Sorghum is the world's fifth

78    biggest crop (after rice, wheat, maize, and barley), belonging to a C4 grass that grows in arid and

79    semi-arid regions (Taylor et al. 2010). Its drought tolerance is a consequence of morphological

80    and anatomical characteristics (i.e., thick leaf wax, deep root system) and physiological

81    responses (i.e., stay-green, osmotic adjustment), is considered as a plant model for drought

82    tolerance in genomic research (Sunita et al. 2011). Recently, the completion of the whole

83    genome assembly of Sorghum (*Sorghum bicolor* L.) makes it possible to identify and analyze

84    cystatin family genes in Sorghum (Paterson et al. 2009). In this study, we aimed to perform a

85    genome-wide identification of *SbCys* family genes in Sorghum and analyze their phylogeny,

86    conserved motifs, structure, *cis*-elements, and expression profile in different tissues. We also

87    explored the expression patterns of *SbCys* genes in response to biotic and abiotic stresses. The

88    results may lay a foundation for further functional analyses of cystatin genes.

89

90    **MATERIALS AND METHODS**

91    **Identification of SbCys family members in Sorghum genome**

92   The identification of SbCys candidates was conducted according to the methods of Lozano et al.

93   (2015) with some modification. The cystatin sequences of *Arabidopsis*, rice, and barely were

94   downloaded from TAIR (http://www.Arabidopsis.org), the Rice Genome Annotation Project

95   (http://rice.plantbiology.msu.edu/index.shtml), and Ensembl database (http://plants.ensembl.org),

96   respectively. The whole-genome sequence of Sorghum was downloaded from Ensembl database

97   (http://plants.ensembl.org). Then predicted proteins from Sorghum genome were scanned using

98   HMMER v3 (http://hmmer.org/) using the Hidden Markov Model (HMM) profile of cystatin

99   (PF00031) from Pfam protein family database (http://pfam.xfam.org/) (Finn et al. 2011). From

100   the proteins obtained using the raw cystatin HMM, a high-quality protein set with a cut-off *e*-

101   value $< 1 \times 10^{-10}$ was aligned and used to construct a Sorghum specific cystatin HMM using

102   hmmbuild from the HMMER v3 suite. Then all proteins with *e*-value < 0.01 were selected by the

103   new Sorghum specific HMM. Cystatin sequences were further filtered based on the closest

104   homolog from *Arabidopsis*, rice and barely using ClustalW and the UNIREF100 sequence

105   database. Proteins without typical domain (Aspartic acid proteinase inhibitor) and reactive site

106   motif (QxVxG) were removed from posterior analysis.

107   **Sequence alignment, structure analysis, and phylogenetic tree construction**

108   The Multiple Expectation for Motif Elicitation (MEME) program was used to identify conserved

109   motifs shared among SbCys proteins. The parameters of MEME were as follows: maximum

110   number of motifs, 10; optimum width, between 6 and 50; and number of repetitions, any.

111   The three-dimensional structures of Sorghum cystatins were modelled by the automated SWISS-

112   MODEL program (http://swissmodel.expasy.org/interactive) (Peitsch 1996). The known crystal

113   structure of rice oryzacystatin I (OC-I) (Nagata et al. 2000) and SiCYS (Hu et al. 2015) were

114   used to construct the homology-based models. Structure analysis was conducted by the RasMol

115    2.7 program (Sayle and Milner-White 1995).

116    A phylogenetic tree was constructed using MEGA X with the maximum likelihood method

117    according to the Whelan and Goldman + freq. Model. Bootstrap analysis was performed by 1000

118    replicates with the p-distance model. The phylogenetic tree was visualized and optimized in

119    Figtree (http://tree.bio.ed.ac.uk/software/figtree/).

120    **Transcript structures, chromosomal location and gene duplication**

121    The genomic structure of each *SbCys* gene was derived from the alignment of their coding

122    sequence to their corresponding genome full-length sequence. The diagrams of these *SbCys*

123    genes were drawn by the Gene Structure Display Server (GSDS, http://gsds.cbi.-pku.edu.cn/)

124    (Hu et al. 2014). The chromosomal locations of *SbCys* genes were retrieved from the

125    Sorghum_bicolor_NCBIv3 map. The genes were plotted on chromosomes using the Map

126    Gene2chromosome (MG2C, version 2.0) tool (http://mg2c.iask.in/). Gene duplication events of

127    *SbCys* family genes were investigated according to the following two criteria: (1) the alignment

128    covered > 75% of the longer gene, (2) the aligned region had an identity > 75%, (3) located in

129    less than 100 kb single region or separated by less than five genes (Gu et al. 2002). For

130    microsynteny analysis, the CDS sequence of every cystatin from *Arabidopsis*, barley, rice, and

131    Sorghum was used as the query to search against all other cystatins using NCBI_blast software

132    with *e*-value $\leq$ 1e$^{-10}$. The Circos software was used to display the results of collinearity gene

133    pairs (Krzywinski et al. 2009).

134    **Calculation of Ka and Ks**

135    To assess the degree of natural selection on *SbCys* genes, the rate ratio of *Ka* (nonsynonymous

136    substitution rate) to *Ks* (synonymous substitution rate) was calculated using KaKs Calculator 2.0

137    (Zhang et al. 2006). The Ka/Ks ratio > 1, < 1, or = 1 indicates positive, negative, or neutral

138    evolution, respectively (Yadav et al. 2015).

139    **Promoter analysis of *SbCys* genes**

140    To investigate the *cis*-regulatory elements in a promoter region, the upstream sequences (1.5 kb)

141    of the start codon in each *SbCys* gene were scanned in the PlantCARE database

142    (http://bioinformatics.psb.ugent.be/webtools/plantcare/html/)    and    New    PLACE

143    (https://www.dna.affrc.go.jp/PLACE/?action=newplace).

144    **Analysis of interaction networks of the SbCys proteins**

145    The functional interacting network models of SbCys proteins were integrated using the web

146    STRING program (http://string-db.org/) based on an *Arabidopsis* association model; the

147    confidence parameters were set at a 0.40 threshold, the number of interactors was set to five

148    interactors. *Arabidopsis* AtCys proteins were mapped to Sorghum SbCys proteins based on their

149    homologous relationship, and the interaction network of SbCys proteins was drawn by

150    Cytoscape_v3.6.0.

151    **Expression analysis of *SbCys* genes under biotic stresses**

152    The RNA-Seq data used for investigating the expression patterns of *SbCys* genes in various

153    tissues were downloaded from NCBI SRA (Sequence Read Archive) database (ERP024508)

154    (Wang et al. 2018). Root, shoot, and whole organism were collected at 14 days after germination.

155    Embryo, endosperm and pericarp were collected at 20 days after pollination. Pollen samples

156    were collected at booting stage. Inflorescences were collected according to the sizes: 1-5 mm, 5-

157    10 mm, and 1-2 cm. Three biological replicates were performed for each plant tissue. RNA was

158    sequenced using the Illumina HiSeq 2500 system to generate 250 bp pair-end reads.

159    RNA-seq data of biotic stresses were obtained from two experiments. The first experiment

160    measured the transcriptome response of a resistant Sorghum (*Sorghum bicolor* L. Moench)

161    infected with *Bipolaris sorghicola* (DRP 000986) (Yazawa et al. 2013). RNA samples were

162    collected at 0, 12 and 24 hours post-inoculation with one biological replicate. RNA-seq was run

163    using Illumina technology to give 100-base-pair single-end reads on a HiSeq2000 system. The

164    second study measured changes in the transcriptome of Sorghum leaves infested by sugarcane

165    aphid (Tetreault et al. 2019). The RNA-seq data were downloaded from the NCBI SRA database

166    (SRP162227). In this study, two treatments (infested and control) were arranged and two

167    Sorghum genotypes (resistant cultivar RTx2783 and susceptible cultivar BCK60) were used.

168    Leaf samples were collected from treated and control plants at 5, 10 and 15 days post sugarcane

169    aphid infestation. Three biological replicates were performed for all treatment and time

170    combinations. RNA was sequenced using the Illumina Hiseq 2500 platform to generate 100 bp

171    single end reads. The accession numbers and sample information were listed in Table S1. The

172    differential expression of *SbCys* genes were investigated by Hisat2 (http:/kim-lab.org/), Htseq

173    (http://www.htseq.org/), and DESeq2 (R package) based on the RNA-seq data (Wen, 2017). The

174    $p \leq 0.05$ and $|logFC| \geq 1.5$ were set as the cut-off criterion.

175    **Plant materials and treatments**

176    Seed of Sorghum (*Sorghum bicolor* L. cv. Jinza 35) were surface sterilized (15 min in 4%

177    NaClO), washed with distilled water several times, and transferred to moist germination paper

178    for 3 days in an incubator at 25 °C. These seedlings were grown in holes of foam floating plastic

179    containers (30 seedlings per container) with constant aeration in Hoagland solution in a growth

180    room with 14 h/30 °C light and 10 h/22 °C dark regime. The nutrient solution was routinely

181    changed every 3 days. At the three-leaf stage (the juvenile phase (Hashimoto et al. 2019)),

182    abiotic stresses including ABA, salinity, and dehydration treatments were initiated according to

183    the procedures described in previous reports (Dugas et al. 2011; Wang et al. 2012; Yan et al.

184   2017). The plants were transferred quickly to the nutrient solution containing 0.1 mM ABA

185   (dissolved in ethanol), 5 μL ethanol (control for ABA treatment), 250 mM sodium chloride

186   (NaCl), or 15% (W/V) polyethylene glycol (PEG) 6,000. The central part of flag leaves from

187   randomly selected Sorghum plants were harvested respectively at 0, 12 and 24 hours post-

188   treatment per trial, and immediately frozen in liquid nitrogen and stored at -80 ˚C prior to RNA

189   isolation. For each treatment at a given time, three biological replicates were used. The leaf

190   samples of 10 plants came from the same container for one biological replicate. That is, three

191   containers were used for three biological replicates respectively.

192   **RNA extraction and qRT-PCR analysis**

193   Total RNA of 100 mg leaf samples was isolated using the "TaKaRa MiniBEST Plant RNA

194   Extraction" Kit (TaKaRa, Dalian, China) following the manufacturer's instructions. Purity and

195   concentration of RNA samples were evaluated by measuring the $A_{260}/A_{230}$ and $A_{260}/A_{280}$ ratios.

196   In order to digest the genomic DNA, the RNAs were treated with RNase-free DNase I. Reverse

197   transcription was performed according to the kit instruction (Promega, Madison, USA). Primer

198   pairs for qRT-PCR analysis were designed by Primer3Plus program

199   (http://www.bioinformatics.nl), and shown in Table S2. A 20 μl reaction volume containing 0.4

200   μl of each primer (forward and reverse), 2 μl 10-fold diluted cDNA, 7.2 μl of nuclease-free water

201   and 10 μl of GoTaq® qPCR Master Mix (Perfect Real Time; Promega). PCR reaction included

202   one cycle at 95 ˚C for 3 min, followed by 39 cycles of 95 ˚C for 15 s, 60 ˚C for 30s and 72 ˚C for

203   20s. The reactions were conducted using CFX96 Real-Time PCR Detection System (Bio-Rad

204   Laboratories, Inc.). Three independent biological replicates and two technical replicates of each

205   sample were performed. Gene-specific amplification of both reference and *cystatin* genes were

206   standardized by the presence of a single, dominant peak in the qRT-PCR dissociation curve

207   analyses. All data were analyzed by CFX Manager Software (Bio-Rad Laboratories, Inc.). The

208   efficiency range of the qRT-PCR amplifications for all of the genes tested was between 91% and

209   100%. The average target (*SbCys*) cT (threshold cycle) values were normalized to reference (*β-*

210   *actin*) cT values. The fold change between treated sample and control was calculated using the

211   slightly modified $2^{-(\Delta\Delta Ct)}$ method as described by Kebrom et al. (2010). A probability of $p \leq 0.05$

212   was considered to be significant.

213

214   **RESULTS**

215   **Identification and analysis of *SbCys* genes**

216   To extensively identify all of SbCys family members in Sorghum, we constructed a Sorghum-

217   specific HMM for the SbCys domain to scan Sorghum genome, and 22 gene candidates were

218   identified. After removing the repetitive and/or incomplete sequences, the rest of SbCys

219   sequences were submitted to Pfam (http://pfam.xfam.org/) and SMART (http://smart.embl-

220   heidelberg.de/) to confirm the conserved domain. Finally, a total of 18 non-redundant SbCys

221   proteins were identified and were serially renamed from *SbCys1* to *SbCys17* according to their

222   location and order in chromosomes. Gene names, gene IDs, chromosomal locations, amino acid

223   numbers and protein sequences were listed in Table S3. The average length of these SbCys

224   proteins was 148 amino acid residues and the length mainly centered on the range of 105 to 240

225   amino acid residues.

226   Chromosome distribution analysis showed that the number of *SbCys* genes on each chromosome

227   is different (Fig. 1). Chromosome 1 had the greatest number of *SbCys* genes (9 genes), followed

228   by chromosomes 9 and 3 (4 and 3 genes, respectively). Chromosomes 2 and 4 had just one

229   *SbCys* gene, whereas chromosomes 5, 6, 7, 8 and 10 had no *SbCys* genes. Half of *SbCys* genes

230  were distributed on chromosome 1, suggesting that *SbCys* genes may have a chromosomal

231  preference.     proposed reason?

**Gene structure analysis of *SbCys* genes**

233  The analysis of exon-intron structure can provide significant information about the gene function,

234  organization and evolution of multiple gene families (Xu et al. 2012). Schematic structures of

235  *SbCys* genes from Sorghum were obtained using the GSDS program (Fig. 2). Among the *SbCys*

236  genes, more than half (12, 66.7%) were intronless, three genes (*SbCys11*, *SbCys15*, and *SbCys16*)

237  had one intron, two genes (*SbCys14* and *SbCys17*) had two introns, and one gene (*SbCys10*) had

238  three introns. These six *SbCys* genes with one or more introns were clustered into one clade,

239  suggesting the evolutionary event may effect on the gene structure (Altenhoff et al. 2012).

**Sequence alignment, protein motifs analysis, and structural predication of SbCys**

241  Alignments of SbCys sequences were carried out to search for amino acid variants that could

242  lead to differences in their inhibitory capability for cysteine proteases. The results were shown in

243  Fig. 3a. N-terminal and C-terminal extensions with varying lengths that presented in several

244  SbCys proteins were not displayed in the comparison. These predicted structures shared many

245  identical residues including α-helix and the four β-sheets (β2-5) (Fig. 3a). Analysis of conserved

246  motifs of SbCys proteins also revealed that some typical conserved motifs could be detected in

247  most SbCys proteins, such as motif 1, 2, 3, and 4, form a fundamental structural combination

248  (Fig. 3b and 3c). Motif 1 was conserved in the central loop region with a consensus sequence of

249  "QxVxG" and could be detected in most SbCys proteins, which played an important role in the

250  inhibitory capacity of cystatins towards their target cysteine proteases (Meriem et al. 2010).

251  Motif 2 contained a particular consensus sequence ([LVI][GA][RQG][WF]AV) that conformed

252  to a predicted secondary α-helix structure (Martinez et al. 2009). The other two typical motifs for

253  SbCys proteins, motif 3 (V[WY][EVG]KPW) and motif 4 ([RK]xLxxF), were firstly described

254  in tobacco (Zhao et al. 2014), were also detected in most SbCys proteins, indicating their

255  conserved and common role in both dicots and monocots. Motif 5 existed only in 3 SbCys family

256  members (SbCys5, SbCys8, and SbCys15,). Details of the 5 conserved motifs were shown in Fig.

257  S1.

258  The predicted three-dimensional structures of the Sorghum cystatins were established using the

259  SWISS-MODEL program based on the known crystal structure of OC-I and SiCYS (Fig. 4).

260  Although these structures were predicted with variable degrees of accuracy, all of Sorghum

261  cystatins shared similar protein structure with rice OC-I (Fig. 4a), excepting SbCys10 that shared

262  similar protein structure with SiCYS (Fig. 4b). In additon, SbCys14 showed a significant

263  variation in its predicted three-dimensional structures, an extra α-helix occurred in the C-terminal

264  region, which probably due to the cystatin contained a C-terminal extension. Two important

265  motifs (the conserve QxVxG motif and W residue) of Sorghum cystatins involved in the

266  interaction with the target cysteine enzymes were also shown in Fig. 4. The predicted structure of

267  SbCys13 showed some distortions in the region of the β2 sheet, probably due to the insertion of a

268  methionine in the first position of the conserved QxVxG motif.

269  **Phylogenetic analysis of *SbCys* genes**

270  Cystatin gene family is highly conserved in both monocots and dicotyledons (Martinez and Diaz,

271  2008). To investigate the phylogenic relationships of SbCys proteins to other known plant

272  cystatins, a multiple sequence alignment of SbCys sequences to the sequences from *Arabidopsis*,

273  rice, and barley was conducted by the ClustalW program. As showed in Fig. 5, these cystatins

274  were categorized into three groups, including Group I, Group II, and Group III. A total of 21

275  cystatins were classified to Group I and 6 cystatins from Sorghum. Group II contained 7

276    cystatins, only one cystatin from Sorghum. The remaining 21 proteins were assigned to Group III

277    and 11 SbCys proteins fell into this group. In addition, some bootstrap values in the phylogenetic

278    tree were low, suggesting that high sequence differentiation in these cystatins occurred.

279    Microsynteny analysis indicated that one orthologous gene pair was identified in the cross of

280    barley and Sorghum, rice and Sorghum, respectively, while no orthologous gene pair between

281    *Arabidopsis* and Sorghum was found (Fig. S2). These data indicated that *SbCys* genes were more

282    closely to those of rice and barley than that of *Arabidopsis*. Interestingly, a pair of *SbCys* genes

283    (*SbCys2-1* and *SbCys2-2*) was confirmed to be tandem duplication in Sorghum (Fig. S2).

284    Analysis of duplicated *SbCys* genes showed that the *Ka*/*Ks* ratios far less than 1, varying from

285    0.0976 to 0.5679 (Table S4), indicating that negative selection occurred in the duplication event.

286    **Promoter analysis of *SbCys* genes**

287    In order to obtain useful information on the regulatory mechanism of cystatin gene expression,

288    the 1.5 kb upstream sequences from the translation start sites of *SbCys* genes were submitted into

289    PlantCARE database to detect the *cis*-elements. Various putative plant regulatory elements in the

290    promoter region of *SbCys* genes were shown in Fig. 6 and Table S5. Several potential regulatory

291    elements involved in stress-related transcription factor-binding sites were found, including G-

292    box, W-box, TC-rich repeats, MBS, heat shock elements (HSEs), and ABA-response element

293    (ABRE). The identified *SbCys* genes possessed at least 1 stress-response-related *cis*-element,

294    suggesting that the expressions of *SbCys* genes were related to these abiotic stresses. All of

295    *SbCys* genes had one or more G-box with the exception of *SbCys9*, implying that these *SbCys*

296    genes could be induced by light stress. 14 *SbCys* genes possessed MBS element, ABRE element

297    was found in 12 *SbCys* genes, and HSE element was located in 10 *SbCys* genes. TC-rich repeats

298    and W-boxes were located in 8 genes, respectively. In addition, Skn-1_motif was conserved in

299    the promoter regions of most *SbCys* genes, indicating these genes were associated with the

300    regulation of seed storage protein gene expression (Strömvik and Fauteux, 2009). The high

301    diversity of the *cis*-acting elements suggested that these *SbCys* genes might have a wide range of

302    functional roles and could be involved in multiple stress responses and growth and development

303    progress (Zhang et al. 2008).

304    **Protein interaction network of SbCys proteins**

305    In this study, the interactions of the SbCys proteins were investigated in an *Arabidopsis*

306    association model using STRING software. As shown in Fig. 7, the interaction network of

307    cystatins showed a complex functional relationship. AtCys2 (corresponding to SbCys12)

308    interacted with stress related proteins (AT1G56280, AT3G19580, AT5G67450, and AtCys1) and

309    growth and development related proteins (AT1G63100 and AT5G04340), AtCys1

310    (corresponding to SbCys11, 15, 16, and 17) interacted with some vacuolar-processing enzyme

311    which involved in processing of vacuolar seed protein precursors into the mature forms, and

312    AtCys5 (corresponding to SbCys1, 2-1, 3, 4, 5, 6, 7, 8, 9, and 13) interacted with several lipid-

313    transfer proteins (AT1G07747, AT1G52415, AT2G16592, AT3G29152, and AT4G12825). The

314    results suggested that cystatins might be associated with many biological processes by protein

315    interactions, such as pollen development, stress responses, and seed maturation (Wang et al.

316    2012).

317    **Expression profile of *SbCys* genes in different Sorghum tissues**

318    To obtain the spatial and temporal expression patterns of all *SbCys* genes, RNA-seq data

319    (ERP024508) were downloaded to explore the expression levels of *SbCys* genes in different

320    tissues including root, stem, whole organism, pollen, endosperm, embryo, inflorescence (1-5mm,

321    1-10mm, and 1-2cm), and pericarp. As shown in Fig. 8 and S3, most *SbCys* genes were

expressed in one tissue at least, except for *SbCys13*, which were barely expressed in any tissue.

The expression patterns of *SbCys* genes in reproductive tissues were significantly difference from vegetable tissues. Such as *SbCys2-1*, *SbCys3*, *SbCys4*, *SbCys5*, *SbCys7*, *SbCys9*, *SbCys12*, and *SbCys17* showed relatively higher expression levels in reproductive tissues including pollen, endosperm, embryo, and pericarp than in vegetable tissues, while the expression of *SbCys7* and *SbCys15* were higher in vegetable tissues than in reproductive tissues. It was worth noting that majority of *SbCys* genes had lower expression levels during inflorescence development excepting for *SbCys17* which displayed higher expression pattern.

**Expression of *SbCys* genes under biotic stresses**

To gain insight into the potential roles of *SbCys* genes in response to *Bipolaris sorghicola* infection and sugarcane aphid infestation, their relative expression patterns were investigated by using the public transcription data from NCBI SRA database (DRP000986 and SRP162227, respectively). As shown in Fig. 9 and 10, the expression patterns of *SbCys* genes were different under the two biotic stresses. In response to *Bipolaris sorghicola* infection, seven *SbCys* genes were induced and only 2 genes (*SbCys12* and *SbCys13*) were suppressed in infected Sorghum leaves compared with control (Fig. 9a). However, under aphid infestation, four *SbCys* genes (*SbCys4*, *SbCys10*, *SbCys11*, and *SbCys14*) were up-regulated and 3 genes (*SbCys1*, *SbCys3*, and *SbCys17*) were down-regulated relative to control in susceptible Sorghum line (BCK60). In resistant Sorghum line (RTx2783), only two *SbCys* genes (*SbCys4* and *SbCys11*) were induced, and the rest were barely expressed in Sorghum leaves with aphid infection (Fig. 9b and 10). These results might suggest that *SbCys* genes played different roles in responding to pathogen infection and aphid infestation.

**Expression profiling of *SbCys* genes under abiotic stresses**

345 We also investigated the expression of *SbCys* genes in response to various abiotic stresses

346 including dehydration, salt shock, and ABA (Fig. 11). Under dehydration stress, seven *SbCys*

347 genes (*SbCys4*, *SbCys5*, *SbCys6*, *SbCys9*, *SbCys10*, *SbCys11*, and *SbCys17*) were induced to

348 present a significant up-regulation from 0 to 24 h, while the expressions of *SbCys2-1*, *SbCys12*,

349 *SbCys15*, and *SbCys16* were decreased. Furthermore, the expressions of 4 *SbCys* genes (*SbCys1*,

350 *SbCys3*, *SbCys8*, and *SbCys14*) displayed an up-down trend from 0 h to 24 h (Fig. 11a). With salt

351 shock treatment, the expressions of *SbCys2-1*, *SbCys3*, *SbCys4*, *SbCys8*, *SbCys10*, and *SbCys11*

352 were significantly up-regulated at all treatment time points, whereas *SbCys16* showed a

353 significant down-regulated trend (Fig. 11b). In addition, *SbCys6*, *SbCys13 SbCys14*, *SbCys15*,

354 and *SbCys17* showed up and down expression trends, but *SbCys5* displayed down and up

355 expression pattern (Fig. 11b). After exogenous ABA treatment, the expressions of 4 *SbCys* genes

356 (*SbCys2-2*, *SbCys3*, *SbCys4*, and *SbCys7*) were significantly up-regulated at all time points, but 9

357 genes (*SbCys1*, *SbCys2-1*, *SbCys5*, *SbCys8*, *SbCys10*, *SbCys11*, *SbCys13*, *SbCys14*, and *SbCys17*)

358 were down-regulated. Additionally, *SbCys12*, *SbCys15*, and *SbCys16* displayed an up-down

359 expression trends (Fig. 11c). Interestingly, all *SbCys* genes were up-regulated in response to one

360 or two stresses except *SbCys4* that was significantly induced under dehydration, salt shock and

361 ABA stresses, suggesting that SbCys4 might play an important role in response to different stress

362 responses.

363

364 **DISCUSSION**

365 Plant cystatins are a group of intrinsic small proteins, whose members play important roles in

366 diverse biological processes and stress responses (Martinez et al. 2016; Meriem et al. 2010).

367 Recently, a large number of sequence data from different plant species have been uploaded in

368    GenBank, which provide convenience for us to describe their characteristics, and several

369    cystatins families have been identified from plants, such as rice, soybean and wheat (Wang et al.

370    2015; Dutt et al. 2016; Yuan et al. 2016). However, little is known about cystatin family in

371    Sorghum. In the present study, we identified 18 *SbCys* genes from Sorghum genome. The

372    number was less than that of *B. distachyon* genome, where 25 *BdCys* members were identified

373    (Subburaj et al. 2017). The 18 members in Sorghum was a larger number than found in rice (11

374    genes) and *Arabidopsis* (7 genes) (Wang et al. 2015), but was similar to soybean (20 members)

375    (Yuan et al. 2016). The difference on the cystatin number might reflect the adaptation of plants

376    to environment.

377    The identified *SbCys* genes were unevenly distributed on chromosomes 1, 2, 3, 4, and 9, and half

378    of them were distributed on chromosome 1 (Fig. 1), suggesting that *SbCys* genes had a

379    chromosomal preference. The uneven distribution of *cystatin* genes in chromosomes was also

380    found in *B. distachyon* genome that the highest number of *BdCys* genes located in chromosome 1

381    (Subburaj et al. 2017). The phnomenon of chromosomal preference was also observed in *Oryza*

382    *sativa* genome, but most *OsCys* genes were dispersed over two chromosomes, chromosome 1

383    and 3 (Wang et al. 2015). Furthermore, several tandem duplication events occurred at

384    chromosomes 1 of *B. distachyon* genome (Subburaj et al. 2017). Two tandem duplication events

385    (*OsCys4*/*OsCys5* and *OsCys6*/*OsCys7*) were found among *OsCys* genes, and existed in

386    chromosomes 1 and 3 (Wang et al. 2015). One tandem duplication event (*SbCys2-1*/*SbCys2-2*)

387    occurred among *SbCys* genes at chromosome 1 (Fig. S2). The tandem duplication events might

388    cause the distinct distribution patterns of *cystatin* genes on the chromosomes (Li et al. 2017).

389    Eighteen *SbCys* genes were divided into three groups based on phylogenetic analysis (Fig. 5).

390    Some conserved motifs among SbCys proteins had been identified by the alignment of the amino

391   acid sequences (Fig. 3). However, the conservation was accompanied with differences in some

392   important amino acids, indicating that SbCys family members might undergo a complex

393   evolutionary history, which would have a significant influence on their respective functions

394   (Abraham et al. 2006). For example, QxVxG motif, could directly enter and interact with the

395   active site of targeted enzymes, were conserved in all SbCys proteins with the exceptions of 5

396   cystatins (SbCys1, SbCys6, SbCys8, SbCys9, and SbCys13) that were partially modified by the

397   insertion or variation in important residues (Fig. 3a). Furthermore, three SbCys proteins (SbCys8,

398   SbCys9, and SbCys13) showed significant variations with other Sorghum cystatins in their

399   predicted three-dimensional structures (Fig. 4). The variations in vital amino acid residues might

400   result in the change in cystatin inhibitory action (Melo et al. 2003). In addition, two novel motifs,

401   motif 3 (V[WY][EVG]KPW) and motif 4 ([RK]xLxxF), firstly described in tobacco (Zhao et al.

402   2014), were also identified in the C-terminalin of many SbCys proteins. The contribution of the

403   two new motifs to cystatin inhibitory action needs to be further studied.

404   During past decades, plant cystatins were reported to play essential roles in inhibiting

405   endogenous and exogenous cysteine proteases activities during seed development (Gaddour et al.

406   2001; Kiyosaki et al. 2007). In the present study, as revealed by RNA-seq data analysis (Fig. 8

407   and S3), the expression levels of several *SbCys* family genes were higher in reproductive tissues

408   than in vegetable tissues, which were consistent with the reports that most cystatins were

409   specifically expressed in developing seeds and played a role in seed development (Dutt et al.

410   2010; Zhao et al. 2014). Moreover, promoter analysis showed that the highly expressed *SbCys*

411   genes in reproductive tissues possessed endosperm expression related *cis*-elements (Skn-1 and

412   GCN4_motif) (Fig. 6 and Table S5). Our protein interaction prediction results also showed that

413   several SbCys proteins could interact with many functional proteins (Fig. 7), implying these

414  cystatins were involved in regulating the gene expression of cereal grain storage proteins (Diaz-

415  Mendoza et al. 2016).

416  Plant cystatins are involved in various biotic stress responses and probably act as defense

417  proteins against pests and pathogen infection (Meriem et al. 2010). At present, some cystatins

418  with insecticidal activity have been isolated from barley, corn, tomato and papaya etc. (Alvarez-

419  Alfageme et al. 2007; Goulet et al. 2008; Kiggundu et al. 2010), and several cystatins having

420  antifungal activities were also isolated from taro, cacao, and wheat (Christova et al. 2006;

421  Pirovani et al. 2010; Chen et al. 2014). Although studies on insecticidal and antifungal activity of

422  plant cystatins have been well established in vitro, the knowledge about their roles in plants in

423  response to biotic stresses is limited. To explore the properties of *SbCys* genes responding to pest

424  and pathogen infection, we conducted the analysis on the expression patterns of *SbCys* genes.

425  The results showed that the expressions of most *cystatin* genes were induced during *Bipolaris*

426  *sorghicola* infection, suggesting these cystatins played functions in inhibiting exogenous

427  cysteine proteases secreted by pathogens to infect plant tissues (Fig. 9a). Interestingly, for

428  sugarcane arthropods infestation, only two genes (*SbCys4* and *SbCys11*) were up-regulated

429  significantly in susceptible and resistant Sorghum lines (Fig. 9b and 10), the expressions of the

430  rest genes were no obvious change or were down-regulated. These differential expression

431  patterns between *SbCys* genes might suggest that some of them had evolved to inhibit specific

432  cysteine proteinases. The exact roles of these *SbCys* genes in insecticidal and antifungal activity

433  in vivo are worthy to be explored in the further study.

434  Another characteristic of cystatin genes is that they are involved in various abiotic stress

435  responses in different plant species, such as rice, barley, and maize (Gaddour et al. 2001;

436  Massonneau et al. 2005; Huang et al. 2007). In *Arabidopsis*, the expression levels of *AtCYS1* and

437     *AtCYS2* were enhanced by high temperature and wounding stresses (Hwang et al. 2010). *AtCYSa*

438     and *AtCYSb* were also induced by different abiotic stresses such as salt, drought, oxidation and

439     cold stresses (Zhang et al. 2008). Velasco-Arroyo et al. (2018) reported that the silence of barley

440     *HvCPI-2* and *HvCPI-4* specifically modified leaf responses to drought stress. Wang et al. (2015)

441     observed the significant change in the expression levels of several rice *OsCYS* genes under cold,

442     drought, salt, and hormone treatments. In the present study, most *SbCys* genes were found to

443     have positive or negative responses to dehydration, salt shock, and ABA stresses. Moreover, the

444     interaction results showed that most cystatins could interact with stresses-related proteins,

445     implying that the cystatins played critical roles in response to diverse stress conditions. Notably,

446     the expression of *SbCys4* was significantly up-regulated under three stress conditions (Fig. 11),

447     suggesting a specific role of SbCys4 in responding to various stress conditions. Promoter

448     analysis indicated that stress-related *cis*-elements were widespread in the promoter region of

449     these cystatin genes (Table S5), and *SbCys4* possessed plenty of stress-related *cis*-elements,

450     including G-box, ABRE, HSE, MBS and TC-rich repeats. These results provide an effective

451     reference for the functional verification of the *SbCys* family genes under abiotic stresses.

452

453     **CONCLUSIONS**

454     In the current study, we identified 18 *SbCys* family genes in Sorghum genome through a

455     genome-wide survey. The chromosomal localization, conserved protein domain, gene structure,

456     the phylogenetic relationship, as well as the interaction network of these *SbCys* genes was

457     systematically analyzed, revealing special characteristics of *SbCys* family genes in Sorghum. The

458     identified *SbCys* genes displayed an uneven distribution in Sorghum chromosomes. All *SbCys*

459     genes shared similar exon/intron organization and conserved motifs. Phylogenetic analysis

460  suggested that Sorghum cystatins had higher homology with monocotyledon than dicotyledon.

461  The variation of amino acids in Sorghum cystatin critical active sites suggested that they might

462  undergo a complex evolutionary process and possess structural and functional divergence. The

463  expression profile of *SbCys* genes in different tissues indicated that most *SbCys* genes were

464  involved in tissue growth and development. Changes in the expression of *SbCys* genes under

465  biotic and abiotic stresses indicated that many *SbCys* genes played important roles in response to

466  unfavorable growth conditions. It was noting that the expression of *SbCys4* was significantly

467  enhanced under biotic and abiotic stresses, suggesting its unique role in mediating the response

468  of Sorghum to adverse environment conditions.

469

470  **REFERENCES**

471  **Abraham Z, Martinez M, Carbonero P, Diaz I. 2006.** Structural and functional diversity

472  within the cystatin gene family of *Hordeum vulgare*. *Journal of Experimental Botany*

473  **57(15):**4245-4255 DOI 10.1093/jxb/erl200.

474  **Altenhoff AM, Studer RA, Robinsonrechavi M, Dessimoz C. 2012.** Resolving the ortholog

475  conjecture: orthologs tend to be weakly, but significantly, more similar in function than

476  paralogs. *PLoS Computational Biology* **8(5):**e1002514 DOI

477  10.1371/journal.pcbi.1002514.

478  **Alvarez-Alfageme F, Martinez M, Pascual-Ruiz S, Castanera P, Diaz I, Ortego F. 2007.**

479  Effects of potato plants expressing a barley cystatin on the predatory bug *Podisus*

480  *maculiventris* via herbivorous prey feeding on the plant. *Transgenic Research* **16:**1-13 DOI

481  10.1007/s11248-006-9022-6.

482  **Belenghi B, Acconcia F, Trovato M, Perazzolli M, Bocedi A, Polticelli F, Ascenzi P,**

483      **Delledonne M. 2010.** AtCYS1, a cystatin from *Arabidopsis thaliana*, suppresses

484      hypersensitive cell death. *European Journal of Biochemistry* **270(12):**2593-604 DOI

485      10.1046/j.1432-1033.2003.03630.x.

486    **Blanca VA, Mercedes DM, Andrea GS, Santamaria B, Estrella M, Miguel TB, Kumlehn G,**

487      **Martinez J, Diaz I. 2018.** Silencing barley cystatins *HvCPI-2* and *HvCPI-4* specifically

488      modifies leaf responses to drought stress. *Plant Cell Environment* **41:**1776-1790 DOI

489      10.1111/pce.13178.

490    **Chen PJ, Senthilkumar R, Jane WN, He Y, Tian Z, Yeh KW. 2014.** Transplastomic

491      *Nicotiana benthamiana* plants expressing multiple defence genes encoding protease

492      inhibitors and chitinase display broad-spectrum resistance against insects, pathogens and

493      abiotic stresses. *Plant Biotechnology Journal* **12(4):**1-13 DOI 10.1111/pbi.12157.

494    **Christova PK, Christov NK, Imai R. 2006.** A cold inducible multidomain cystatin from winter

495      wheat inhibits growth of snow mold fungus, *Microdochium nivale*. *Planta* **223:**1207-1218

496      DOI 10.1007/s00425-005-0169-9.

497    **Christova PK, Christov NK, Mladenov PV, Imai R. 2018.** The wheat multidomain cystatin

498      TaMDC1 displays antifungal, antibacterial, and insecticidal activities in planta. *Plant Cell*

499      *Reports* **37:**923-932 DOI 10.1007/s00299-018-2279-4.

500    **Diaz-Mendoza M, Dominguez-Figueroa JD, Velasco-Arroyo B, Cambra I, Gonzalez-**

501      **Melendi P, Lopez-Gonzalvez A, Garcia A, Hensel G, Kumlehn J, Diaz I, Martinez**

502      **M. 2016.** HvPap-1 C1A protease and HvCPI-2 cystatin contribute to barley grain filling

503      and germination. *Plant Physiology* **170:**2511-2524. DOI 10.1104/pp.15.01944.

504    **Díazmendoza M, Velascoarroyo B, Gonzálezmelendi P, Martínez M, Díaz I. 2014.** C1A

505      cysteine protease-cystatin interactions in leaf senescence. *Journal of Experimental*

506     *Botany* **65(14):**3825-33 DOI 10.1093/jxb/eru043.

507     **Dugas DV, Monaco MK, Olson A, Klein RR, Kumari S, Ware D, Klein PE. 2011.** Functional

508         annotation of the transcriptome of *Sorghum bicolor* in response to osmotic stress and

509         abscisic acid. *BMC Genomics* **12:**514 DOI 10.1186/1471-2164-12-514.

510     **Dutt S, Singh VK, Marla SS, Kumar A. 2010.** In silico analysis of sequential, structural and

511         functional diversity of wheat cystatins and its implication in plant defense. *Genomics*

512         *Proteomics Bioinformatics* **8(1):**42-56 DOI 10.1016/S1672-0229(10)60005-8.

513     **Finn RD, Clements J, Eddy SR. 2011.** HMMER web server: interactive sequence similarity

514         searching. *Nucleic Acids Research* **39:**29-37 DOI 10.1093/nar/gkr367.

515     **Gaddour K, Carbajosa JV, Lara P, Almoneda PI, Diaz I, Carbonero P. 2001.** A constitutive

516         cystatin-encoding gene from barley (Icy) responds differentially to abiotic stimuli. *Plant*

517         *Molecular Biology* **45:**599-608 DOI 10.1023/a:1010697204686.

518     **Goulet MC, Dallaire C, Vaillancourt LP, Khalf M, Badri AM, Preradov A, Duceppe MO,**

519         **Cloutier GC, Michaud CD. 2008.** Tailoring the specificity of a plant cystatin toward

520         herbivorous insect digestive cysteine proteases by single mutations at positively selected

521         amino acid sites. *Plant Physiology* **146:**1010-1019 DOI 10.2307/40065908.

522     **Gu Z, Cavalcanti A, Chen FC, Bouman P, Li WH. 2002.** Extent of gene duplication in the

523         genomes of *Drosophila*, nematode, and yeast. *Molecular Biology Evolution* **19(3):**256-

524         262 DOI 10.1093/oxfordjournals.molbev.a004079.

525     **Hashimoto S, Tezuka T, Yokoi S. 2019.** Morphological changes during juvenile−to−adult phase

526         transition in Sorghum. *Planta* **250:**1557-1566 DOI 10.1007/s00425-013-1895-z.

527     **Hu B, Jin J, Guo AY, Zhang H, Luo J, Gao G. 2014.** GSDS 2.0: an upgraded gene feature

528         visualization server. *Bioinformatics* **31(8):**1296 DOI 10.1093/bioinformatics/btu817.

529     **Hu YJ, Irene D, Lo CJ, Cai YL, Tzen TC, Lin TH, Chyan CL. 2015.** Resonance assignments

530          and secondary structure of a phytocystatin from *Sesamum indicum*. *Biomolecular NMR*

531          *Assignments* **9:**309-311 DOI 10.1007/s12104-015-9598-y.

532     **Huang Y, Xiao B, Xiong L. 2007.** Characterization of a stress responsive proteinase inhibitor

533          gene with positive e.ect in improving drought resistance in rice. *Planta* **226:**73-85 DOI

534          10.2307/23389651.

535     **Hwang JE, Hong JK, Lim CJ, Chen H, Je J, Yang KA, Kim DY, Choi YJ, Lee SY, Lim CO.**

536          **2010.** Distinct expression patterns of two *Arabidopsis* phytocystatin genes, AtCYS1 and

537          AtCYS2, during development and abiotic stresses. *Plant Cell Reports* **29:**905-915 DOI

538          10.1007/s00299-010-0876-y.

539     **Jenko S, Dolenc I, Guncar G, Dobersek A, Podobnik M, Turk D. 2003.** Crystal structure of

540          Stefin A in complex with cathepsin H: N-terminal residues of inhibitors can adapt to the

541          active sites of endo- and exopeptidases. *Journal Molecular Biology* **326(3):**875-885 DOI

542          10.1016/s0022-2836(02)01432-8.

543     **Kebrom TH, Brutnell TP, Finlayson SA. 2010.** Suppression of sorghum axillary bud

544          outgrowth by shade, phyB and defoliation signalling pathways. *Plant Cell Environment*

545          **33(1):**48-58 DOI 10.4161/psb.5.3.11186.

546     **Kiggundu A, Muchwezi J, Van C, Viljoen A, Vorster J, Schlüter U, Kunert K, Michaud D.**

547          **2010.** Deleterious effects of plant cystatins against the banana weevil *Cosmopolites*

548          *sordidus*. *Arch Insect Biochemistry Physiology* **73(2):**87-105 DOI 10.1002/arch.20342.

549     **Kiyosaki T, Matsumoto I, Asakura T, Funaki J, Kuroda M, Misaka T, Arai S, Abe K. 2007.**

550          Gliadain, a gibberellin-inducible cysteine proteinase occurring in germinating seeds of

551          wheat, *Triticum aestivum* L., specifically digests gliadin and is regulated by intrinsic

552        cystatins. *FEBS Journal* **164:**470-477 DOI 10.1111/j.1742-4658.2007.05749.x.

553    **Krzywinski M, Schein J, Birol I, Connors J, Gascoyne R, Horsman D, Jones SJ, Marra MA.**

554        **2009.** Circos: an information aesthetic for comparative genomics. *Genome research*

555        **19(9):**1639-1645 DOI 10.1101/gr.092759.109.

556    **Li J, Yang XW, Li YC, Niu JS, He DX. 2017.** Proteomic analysis of developing wheat grains

557        infected by powdery mildew (*Blumeria graminis* f.sp. *tritici*). *Journal of Plant*

558        *Physiology* **215:**140-153 DOI 10.1016/j.jplph.2017.06.003.

559    **Li SF, Su T, Cheng GQ, Wang BX, Li X, Deng CL, Gao WJ. 2017.** Chromosome evolution in

560        connection with repetitive sequences and epigenetics in plants. *Genes* **8:**290 DOI

561        10.3390/genes8100290.

562    **Lima AM, dos Reis SP, de Souza CR. 2015.** Phytocystatins and their potential to control plant

563        diseases caused by fungi. *Protein and Peptied Letters* **22:**104-111 DOI

564        10.2174/0929866521666140418101711.

565    **Lozano R, Hamblin MT, Prochnik S, Jannink JL. 2015.** Identification and distribution of the

566        NBS-LRR gene family in the Cassava genome. *BMC Genomics* **16(1):**360 DOI

567        10.1186/s12864-015-1554-9.

568    **Margis R, Reis EM, Villeret V. 1998.** Structural and phylogenetic relationships among plant

569        and animal cystatins. *Arch Biochemistry Biophysics* **359(1):**24-30 DOI

570        10.1006/abbi.1998.0875.

571    **Martinez M, Abraham Z, Gambardella M, Echaide M, Carbonero P, Diaz I. 2005.** The

572        strawberry gene Cyf1 encodes a phytocystatins with antifungal activity. *Journal of*

573        *Experimental Botany* **56:**1821-1829 DOI 10.1093/jxb/eri172.

574    **Martinez M, Cambra I, Carrillo L, Diazmendoza M, Diaz I. 2009.** Characterization of the

575        entire cystatin gene family in barley and their target cathepsin L-like cysteine-proteases,

576        partners in the hordein mobilization during seed germination. *Plant Physiology*

577        **151(3):**1531-1545 DOI 10.1104/pp.109.146019.

578        **Martinez M, Diazmendoza M, Carrillo L, Diaz I. 2007.** Carboxy terminal extended

579        phytocystatins are bifunctional inhibitors of papain and legumain cysteine proteinases.

580        *FEBS Letters* **581(16):**2914-2918 DOI 10.1016/j.febslet.2007.05.042.

581        **Martinez M, Diaz I. 2008.** The origin and evolution of plant cystatins and their target cysteine

582        proteinases indicate a complex functional relationship. *BMC Evolutionary Biology*

583        **8(1):**198-210 DOI 10.1186/1471-2148-8-198.

584        **Martinez M, Santamaria ME, Diazmendoza M, Arnaiz A, Carrillo L, Ortego F, Diaz I.**

585        **2016.** Phytocystatins: defense proteins against phytophagous insects and acari.

586        *International Journal of Molecular Sciences* **17(10):**1747-1763 DOI

587        10.3390/ijms17101747.

588        **Massonneau A, Condamine P, Wisniewski J P, Zivy M, Rogowsky PM. 2005.** Maize

589        cystatins respond to developmental cues, cold stress and drought. *Biochimica et Biophysica*

590        *Acta* **1729:**186-199 DOI 10.1016/j.bbaexp.2005.05.004.

591        **Melo FR, Mello MO, Franco OL, Rigden DJ, Mello LV, Genú AM, Silvafilho MC, Gleddie**

592        **S, Grossidesá MF. 2003.** Use of phage display to select novel cystatins specific for

593        *Acanthoscelides obtectus* cysteine proteinases. *BBA-Proteins Proteomics* **1651(1):**146-

594        152 DOI 10.1016/S1570-9639(03)00264-4.

595        **Meriem B, Urte S, Juan V, Marie-Claire G, Dominique M. 2010.** Plant cystatins. *Biochimie*

596        **92(11):**1657-1666 DOI 10.1016/j.biochi.2010.06.006.

597        **Nagata K, Kudo N, Abe K, Arai S, Tanokura M. 2000.** Three-dimensional solution structure

598   of oryzacystatin-I, a cysteine proteinase inhibitor of the rice, *Oryza sativa* L. japonica.

599   *Biochemistry* **39:**14753-14760 DOI 10.1021/bi0006971.

600   **Paterson AH, Bowers JE, Bruggmann R, Dubchak I, Grimwood J, Gundlach H, Haberer G,**

601   **Hellsten U, Mitros T, Poliakov A. 2009.** The *Sorghum bicolor* genome and the

602   diversification of grasses. *Nature* **457(7229):**551-556 DOI 10.1038/nature07723.

603   **Peitsch M. 1996.** ProMod and Swiss-Model: internet-based tools for automated comparative

604   protein modeling. *Biochemical Society Transactions* **224:**274-279 DOI

605   10.1042/bst0240274.

606   **Pfaffl MW. 2001.** A new mathematical model for relative quantification in real-time RT-PCR.

607   *Nucleic Acids Research* **29:**e45 DOI 10.1093/nar/29.9.e45.

608   **Pirovani CP, Santiago AS, Santos LS, Micheli F, Margis R, Silva Gesteira R, Alvim FC,**

609   **Pereira GAG, Mattos JC. 2010.** *Theobroma cacao* cystatins impair *Moniliophthora*

610   *perniciosa* mycelial growth and are involved in postponing cell death symptoms. *Planta*

611   **232(6):**1485-1497 DOI 10.2307/23391912.

612   **Sayle R, Milner-White EJ. 1995.** RasMol: biomolecular graphics for all. *Trends in Biochemical*

613   *Sciences* **20:**374 DOI 10.1016/S0968-0004(00)89080-5.

614   **Song C, Kim T, Chung WS, Lim CO. 2017.** The *Arabidopsis* phytocystatin AtCYS5 enhances

615   seed germination and seedling growth under heat stress conditions. *Molecular Cells*

616   **40(8):**577-586 DOI 10.14348/molcells.2017.0075.

617   **Strömvik MV, Fauteux F. 2009.** Seed storage protein gene promoters contain conserved DNA

618   motifs in *Brassicaceae*, *Fabaceae* and *Poaceae*. *BMC Plant Biology* **9:**126 DOI

619   10.1186/1471-2229-9-126.

620   **Stubbs MT, Laber B, Bode W, Huber R, Jerala R, Lenarcic B, Turk V. 1990.** The refined

621    2.4 A X-ray crystal structure of recombinant human stefin B in complex with the cysteine

622    proteinase papain: a novel type of proteinase inhibitor interaction. *Embo Journal*

623    **9(6):**1939-1947 DOI 10.1002/j.1460-2075.1990.tb08321.x.

624    **Subburaj S, Zhu D, Li X, Hu Y, Yan Y. 2017.** Molecular characterization and expression

625    profiling of *Brachypodium distachyon* L. cystatin genes reveal high evolutionary

626    conservation and functional divergence in response to abiotic stress. *Frontiers in Plant*

627    *Science* **8:**743 DOI 10.3389/fpls.2017.00743.

628    **Sunita K, Klein RR, Andrew O, Monaco MK, Dugas DV, Doreen W, Klein PE. 2011.**

629    Functional annotation of the transcriptome of *Sorghum bicolor* in response to osmotic

630    stress and abscisic acid. *BMC Genomics* **12(1):**514-514 DOI 10.1186/1471-2164-12-514.

631    **Tan Y, Yang Y, Li C, Liang B, Li M, Ma F. 2017.** Overexpression of *MpCYS4*, a phytocystatin

632    gene from *Malus prunifolia* (Willd.) Borkh., delays natural and stress-induced leaf

633    senescence in apple. *Plant Physiology Biochemistry* **115:**219-28 DOI

634    10.1016/j.plaphy.2017.03.025.

635    **Taylor SH, Hulme SP, Rees M, Ripley BS, Woodward FI, Osborne CP. 2010.**

636    Ecophysiological traits in $C_3$ and $C_4$ grasses: A phylogenetically controlled screening

637    experiment. *New Phytologist* **185(3):**780-791 DOI 10.1111/j.1469-8137.2009.03102.x.

638    **Tetreault HM, Grover S, Scully ED, Gries T, Palmer N, Sarath G, Louis J, Sattler SE. 2019.**

639    Global responses of resistant and susceptible Sorghum (*Sorghum bicolor*) to sugarcane

640    aphid (*Melanaphis sacchari*). *Frontiers in Plant Science* **10:**145 DOI

641    10.3389/fpls.2019.00145.

642    **Valdes-Rodriguez S, Galvan-Ramirez JP, Guerrero-Rangel A, Cedro-Tanda A. 2015.**

643    Multifunctional amaranth cystatin inhibits endogenous and digestive insect cysteine

644        endopeptidases: A potential tool to prevent proteolysis and for the control of insect pests.

645        *Biotechnology Applied Biochemistry* **62:**634-641 DOI 10.1002/bab.1313.

646 **Velasco-Arroyo B, Diaz-Mendoza M, Gomez-Sanchez A, Moreno-Garcia B, Santamaria**

647        **ME, Torija-Bonilla M, Hensel G, Kumlehn J, Martinez M, Diaz I 2018.** Silencing

648        barley cystatins HvCPI-2 and HvCPI-4 specifically modifies leaf responses to drought

649        stress. *Plant Cell and Environment* **41(8):**1776-1790 DOI 10.1111/pce.13178.

650 **Wang B, Regulsk M, Tseng E, Olson A, Goodwin S, McCombie WR, Ware D. 2018.** A

651        comparative transcriptional landscape of maize and Sorghum obtained by single-

652        molecule sequencing. *Genome Research* **28(6):**921-928 DOI 10.1101/gr.227462.117.

653 **Wang HW, Hwang SG, Karuppanapandian T, Liu AH, Kim W, Jang CS. 2012.** Insight into

654        the molecular evolution of non-specific lipid transfer proteins via comparative analysis

655        between rice and sorghum. *DNA Research* **19:**179-194 DOI 10.1093/dnares/dss003.

656 **Wang W, Zhao P, Zhou XM, Xiong HX, Sun MX. 2015.** Genome-wide identification and

657        characterization of cystatin family genes in rice (*Oryza sativa* L.). *Plant Cell Reports*

658        **34(9):**1579-1592 DOI 10.1007/s00299-015-1810-0.

659 **Wen G. 2017.** A simple process of RNA-Sequence analyses by Hisat2, Htseq and DESeq2.

660        *International Conference* **Pp:**11-15 DOI 10.1145/3143344.3143354.

661 **Xu G, Guo C, Shan H, Kong H. 2012.** Divergence of duplicate genes in exon-intron structure.

662        *PNAS* **109(4):**1187-1192 DOI 10.1073/pnas.1109047109.

663 **Yadav CB, Bonthala VS, Muthamilarasan M, Pandey G, Khan Y, Prasad M. 2015.**

664        Genome-wide development of transposable elements-based markers in foxtail millet and

665        construction of an integrated database. *DNA Research* **22:**79-90 DOI

666        10.1093/dnares/dsu039.

667 **Yan S, Li SJ, Zhai GW, Lu P, Deng H, Zhu S, Huang RL, Shao JF, Tao YZ, Zou GH. 2017.**

668 Molecular cloning and expression analysis of duplicated polyphenol oxidase genes reveal

669 their functional differentiations in Sorghum. *Plant Science* **263:**23-30 DOI

670 10.1016/j.plantsci.2017.07.002.

671 **Yazawa T, Kawahigashi H, Matsumoto T, Mizuno H. 2013.** Simultaneous transcriptome

672 analysis of Sorghum and *Bipolaris sorghicola* by using RNA-seq in combination with *De*

673 *novo* transcriptome assembly. *PLoS One* **8(4):**e62460 DOI

674 10.1371/journal.pone.0062460.

675 **Yuan S, Li R, Wang L, Chen H, Zhang C, Chen L, Hao Q, Shan Z, Zhang X, Chen S. 2016.**

676 Search for nodulation and nodule development-related cystatin genes in the genome of

677 soybean (*Glycine max*). *Frontiers in Plant Science* **7:**1595 DOI 10.3389/fpls.2016.01595.

678 **Zhang X, Liu S, Takano T. 2008.** Two cysteine proteinase inhibitors from *Arabidopsis thaliana*,

679 AtCYSa and AtCYSb, increasing the salt, drought, oxidation and cold tolerance. *Plant*

680 *Molecular Biology* **68:**131-143 DOI 10.1007/s11103-008-9357-x.

681 **Zhao P, Zhou XM, Zou J, Wang W, Wang L, Peng XB, Sun MX. 2014.** Comprehensive

682 analysis of cystatin family genes suggests their putative functions in sexual reproduction,

683 embryogenesis, and seed formation. *Journal of Experimental Botany* **65(17):**5093-5108

684 DOI 10.1093/jxb/eru274.

685 **Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J. 2006.** KaKs_Calculator 2.0: Calculating

686 Ka and Ks through model selection and model averaging. *Genomics Proteomics*

687 *Bioinformatics* **4(4):**259-263 DOI 10.1016/S1672-0229(10)60008-3.

688

689

690

691

# Figure 1

Chromosome localization of *SbCys* genes.

Chromosome number is indicated at the top of each bar. The size of chromosome was
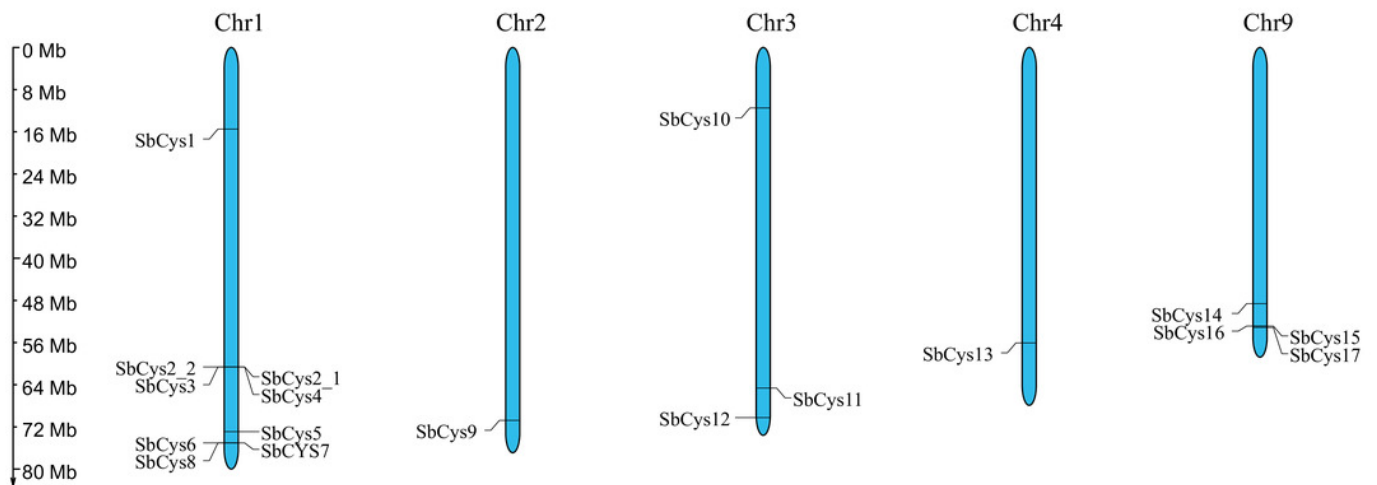labeled on the left of the figure.

# Figure 2

Phylogenetic relationship and gene structure of *SbCys* genes.

A phylogenetic tree was constructed using MEGA X by the maximum likelihood method with 1000 bootstrap replicates. Exon/intron structures were identified by online tool GSDS. Lengths of exons and introns of each *SbCys* genes were exhibited proportionally. Exons and introns are shown by blue bars and black horizontal lines, respectively.

# Figure 3

The amino acid alignment and conserved motifs distribution of SbCys.

(A) The locations of the secondary structures (α-helix and β-sheets) were included. The main cystatin conserved motifs are in black boxes. The strong and weak conservative changes in amino acids are marked by dark gray and light gray font, respectively. (B) The motifs were identified by MEME. Each motif was represented by one color box. (C) Conserved protein motif 1 (QxVxG), motif 2 ( LARFAV and G-residue), motif 3 (W-residue), motif 4 ([RK]xLxxF), and motif 5(P-residue) presented in the variable region of cystatin genes.

# Figure 4

The three-dimensional structure prediction of Sorghum cystatins.

(A) The three-dimensional structures of SbCys proteins were predicted using the automated SWISS-MODEL program with OC-I as a template. (B) The three-dimensional structure of SbCys10 was predicted using the automated SWISS-MODEL program with SiCYS as a template. Two important motifs involved in the interaction with the target enzymes are indicated: the reactive site (asterisks) and W residue (crosses).

# Figure 5

Phylogenetic relationships of the cystatins from *Arabidopsis*, rice, barley and Sorghum.

The phylogenetic tree was constructed by MEGA X with the maximum likelihood method. The numbers at the nodes indicate the bootstrap values. Gene names with black, red, and blue represented Group I, Group II, and Group III, respectively.

# Figure 6

The distribution of *cis*-elements in the 1.5 kb upstream promoter regions of *SbCys* genes.

The *cis*-elements in the promoter region of *SbCys* genes were predicted using PlantCARE database ( http://bioinformatics.psb.ugent.be/webtools/plantcare/html/ ). Different *cis*-elements were represented by different shapes and colors.

# Figure 7

The interaction networks of SbCys proteins according to the orthologs in *Arabidopsis*.

Functional interacting network models were integrated using the STRING tool, and the confidence parameters were set at a 0.40 threshold. Homologous genes in Sorghum and *Arabidopsis* are shown in black and red, respectively.
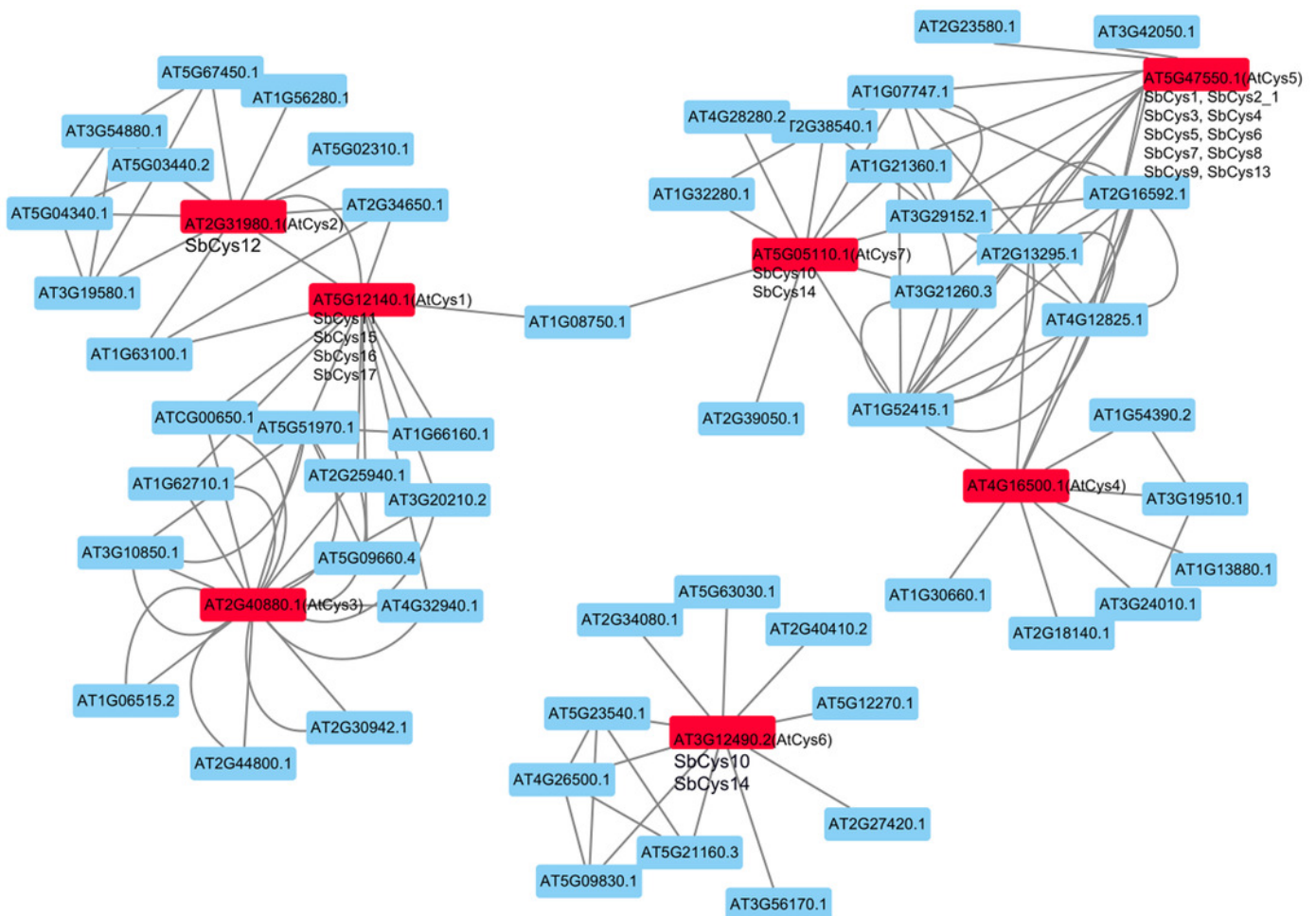
# Figure 8

Hierarchical clustering of the expression profiles of *SbCys* genes in different tissues .

Different tissues are exhibited below each column. Root, shoot, and whole organism belonged to vegetable tissues were collected at 14 days after Sorghum seed germination. Reproductive tissues included embryo , endosperm and pericarp were collected at 20 days after pollination; pollens at booting stage; Inflorescences based on sizes: 1-5 mm, 5-10 mm, and 1-2 cm. Log transform data was used to create the heatmap. The scale bar represented the fold change (color figure online). Blue blocks represented the lower expression level and red blocks represented the higher expression level.
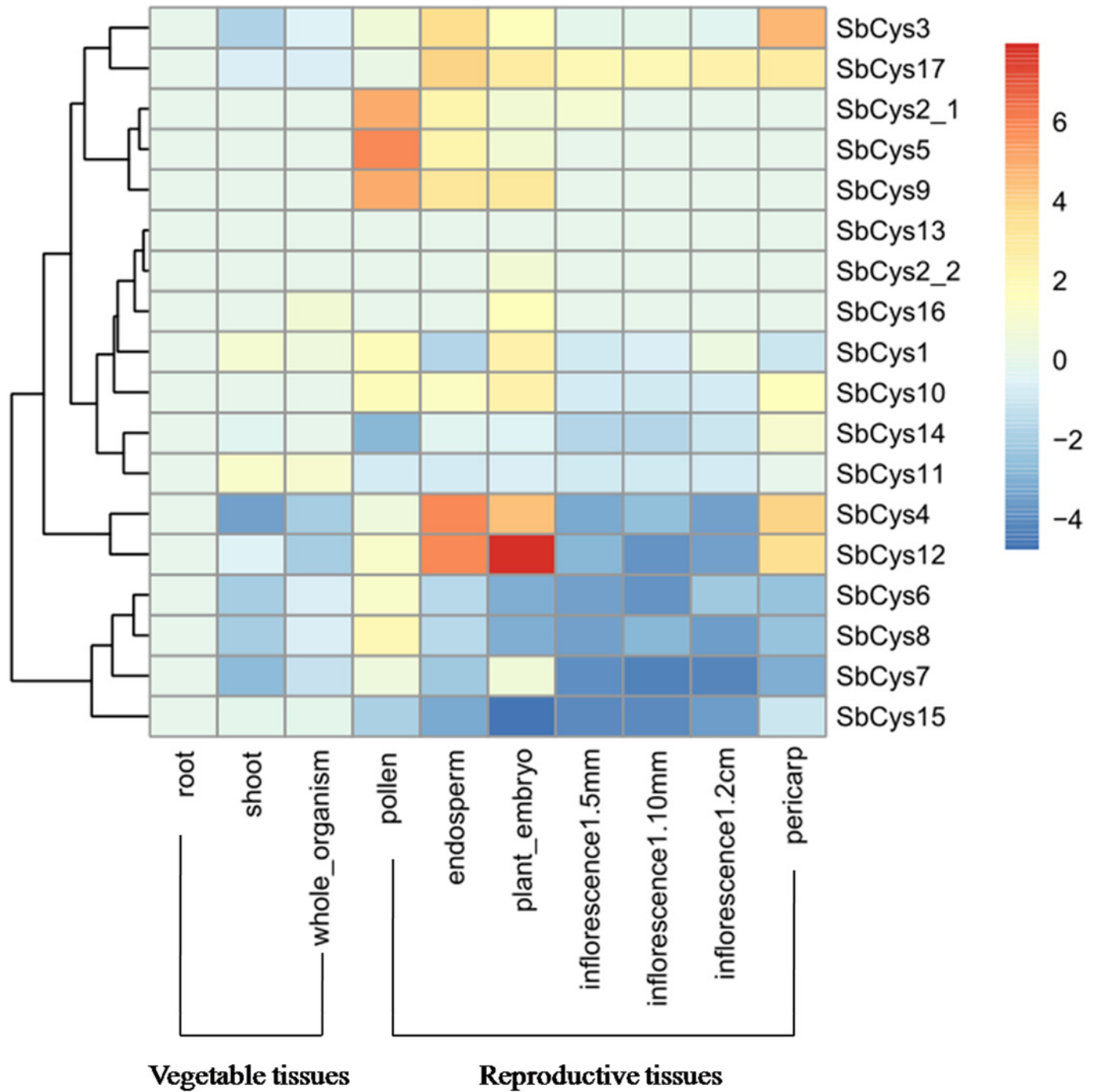
# Figure 9

Hierarchical clustering of the expression profiles of *SbCys* genes under biotic stresses.

(A) The expression changes in *SbCys* genes at 0, 12, and 24 hours with *Bipolaris sorghicola* infection. (B) The expression changes of *SbCys* genes at 5, 10, 15 days with sugarcane aphid infestation. Log transform data was used to create the heatmap. The scale bar represents the fold change (color figure online). Blue blocks indicate low expression and red blocks indicate high expression (color figure online).
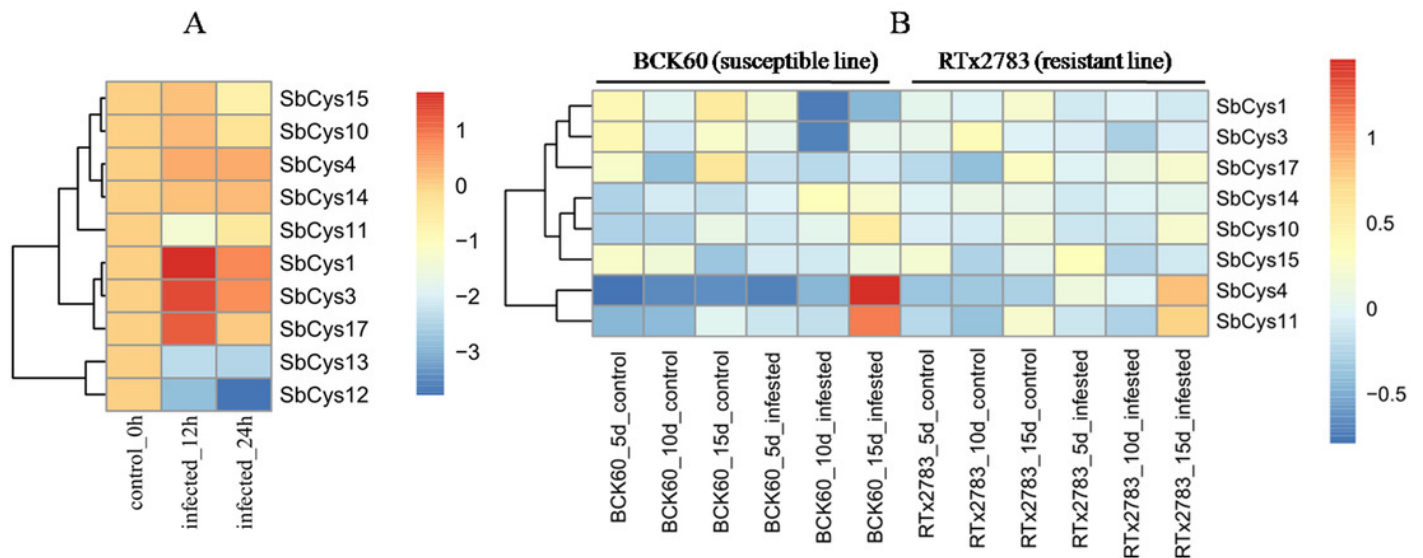
# Figure 10

Expression profiles of *SbCys* genes at 5, 10, and 15 days with sugarcane aphid infection.
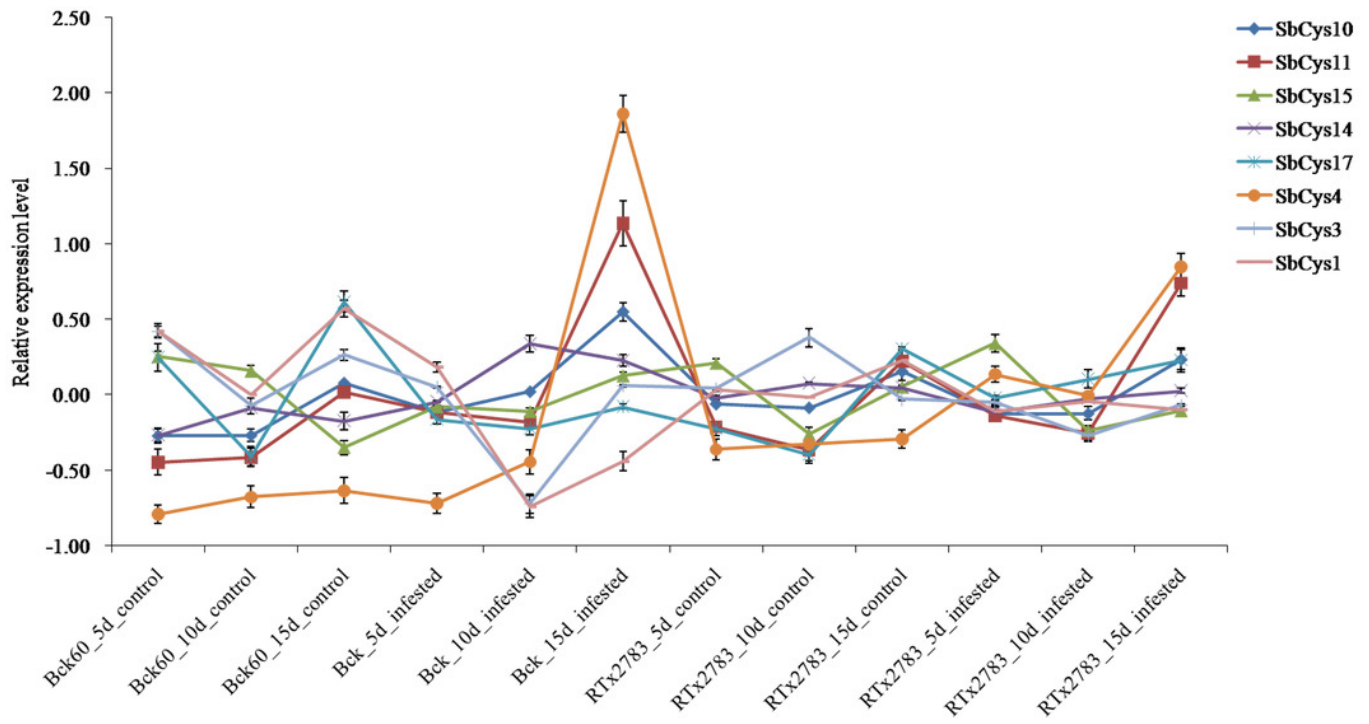
# Figure 11

Expression patterns of *SbCys* genes under (A) dehydration (PEG 6,000) treatment, (B) salt shock (NaCl) treatment, and (C) ABA treatment.

qRT-PCR was used to investigate the expression levels of each *SbCys* gene. To visualize the relative expression levels data, 0 h at each treatment was normalized as "1". * indicated significant differences in comparison with the control at $p \leq 0.05$. ** indicated significant differences in comparison with the control at $p \leq 0.01$.

A. PEG 6,000 treatment

B. NaCl treatment

C. ABA treatment