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

Jie Li [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 Author: Dexian He
Email address: 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 were distributed unevenly on chromosomes, and two genes were involved in a 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 different tissues, and most genes displayed higher expression levels in reproductive tissues than in vegetative tissues, indicating that the *SbCys* genes participated in the regulation of seed formation. Furthermore, the expression profiles of the *SbCys* genes revealed that 7 cystatin family genes were induced during *Bipolaris sorghicola* infection and 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, and ABA stresses). 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 responses of Sorghum to adverse environmental impact. Our results provide new insights into the structural and functional characteristics of the *SbCys* gene family, which lay 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. Moench)

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

**ABSTRACT**

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 were distributed unevenly on chromosomes, and two genes were involved in a 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 different tissues, and most genes displayed higher expression levels in reproductive tissues than in vegetative tissues, indicating that the *SbCys* genes participated in the regulation of seed formation. Furthermore, the expression profiles of the *SbCys* genes revealed that 7 cystatin family genes were induced during *Bipolaris sorghicola* infection and 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, and ABA stresses). 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 responses of Sorghum to adverse environmental impact. Our results provide new insights into the structural and functional characteristics of the *SbCys* gene family, which lay the foundation for better understanding the roles and regulatory mechanism of Sorghum cystatins in seed development and responses to different stress conditions.

**INTRODUCTION**

Cystatins are competitive and reversible inhibitors of cystein proteases from families C1A and

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

47  2014). Cystatins are categorized into three groups, including stefins without disulfide bonds

48  (Group I), cystatins with four conserved Cys residues forming two disulfide bonds (Group II),

49  and kininogens with repeated, stefin-like domains (Group III) (Meriem et al. 2010). Cystatins are

50  widely distributed in animal and plant systems (Tremblay et al. 2019). Based on their primary

51  sequence homology, cystatins contain three signature motifs including a QxVxG reactive site, a

52  tryptophan residue (W) located downstream of the reactive site, and one or two glycine (G)

53  residues in the flexible N terminus of the protein. These three motifs are important for the

54  cystatin inhibitory mechanism (Martinez et al. 2009). In addition, a consensus sequence ([LVI]-

55  [AGT]-[RKE]-[FY]-[AS]-[VI]-x-[EDQV]-[HYFQ]-N) in cystatins is conformed to a predicted

56  secondary α-helix structure (Meriem et al. 2010). Most plant cystatins are small proteins with a

57  molecular mass in the 12- to 16-kD range (Meriem et al. 2010). Some plant cystatins contain a

58  C-terminal extension that raises their molecular weights up to 23 kDa. The longer C-terminal

59  extensions are thought to be involved in the inhibition of cysteine protease activities in the

60  peptidase C13 family (Martinez et al. 2007; Martinez and Diaz, 2008).

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

62  proteases during plant growth and development, senescence, and programmed cell death

63  (Belenghi et al. 2010; Díazmendoza et al. 2014; Zhao et al. 2014). Additionally, plant cystatins

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

65  2016). For example, several publications reported the inhibition of recombinant cystatins on the

66  growth of some pests and fungi (Lima et al. 2015; Tremblay et al. 2019). Tomato plants over-

67  expressing the wheat cystatin *TaMDC1* displayed a broad stress resistance to bacterial pathogen,

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

69 2018). The inhibition of amaranth cystatin on the digestive insect cysteine endopeptidases was

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

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

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

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

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

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

76 *Arabidopsis*, rice, soybean, wheat, *Populus trichocarpa*, and *Brachypodium distachyon*

77 (Martinez and Diaz, 2008; Wang et al. 2015; Yuan et al. 2016; Dutt et al. 2016; Subburaj et al.

78 2017). However, a genome-wide study of cystatin family genes in Sorghum has not yet been

79 performed. Sorghum is the world's fifth biggest crop (after rice, wheat, maize, and barley),

80 belonging to a C4 grass that grows in arid and semi-arid regions (Taylor et al. 2010). Its drought

81 tolerance is a consequence of morphological and anatomical characteristics (i.e., thick leaf wax,

82 deep root system) and physiological responses (i.e., stay-green, osmotic adjustment). Hence,

83 Sorghum is an excellent model plant for the study of plant response to drought stress (Sunita et al.

84 2011). Recently, the completion of the whole genome assembly of Sorghum (*Sorghum bicolor* L.

85 Moench) makes it possible to identify and analyze cystatin family genes in Sorghum (Paterson et

86 al. 2009). In this study, we aimed to perform a genome-wide identification of *SbCys* family

87 genes in Sorghum and analyze their phylogeny, conserved motifs, structure, *cis*-elements, and

88 expression profile in different tissues. We also explored the expression patterns of *SbCys* genes

89 in response to biotic and abiotic stresses. The results may lay a foundation for further functional

90 analyses of cystatin genes.

91

92 **MATERIALS AND METHODS**

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

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

95 (2015) with some modification. The cystatin sequences of *Arabidopsis*, rice (*Oryza sativa*), and

96 barely (*Hordeum vulgare*) were downloaded from TAIR (http://www.Arabidopsis.org), the Rice

97 Genome Annotation Project (http://rice.plantbiology.msu.edu/index.shtml), and Ensembl

98 database (http://plants.ensembl.org/Hordeum_vulgare/Info/Index), respectively. The whole-

99 genome sequence of Sorghum was downloaded from Ensembl database

100 (http://plants.ensembl.org/Sorghum_bicolor/Info/Index). Then predicted proteins from the

101 Sorghum genome were scanned using HMMER v3 (http://hmmer.org/) using the Hidden Markov

102 Model (HMM) profile of cystatin (PF00031) from the Pfam protein family database

103 (http://pfam.xfam.org/) (Finn et al. 2011). From the proteins obtained using the raw cystatin

104 HMM, a high-quality protein set with a cut-off $e$-value $< 1 \times 10^{-10}$ was aligned and used to

105 construct a Sorghum specific cystatin HMM using hmmbuild from the HMMER v3 suite. Then

106 all proteins with $e$-value $< 0.01$ were selected by the new Sorghum specific HMM. Cystatin

107 sequences were further filtered based on the closest homolog from *Arabidopsis*, *Oryza sativa*,

108 and *Hordeum vulgare* using ClustalW and the UNIREF100 sequence database. Proteins that

109 have no typical domain (Aspartic acid proteinase inhibitor) and reactive site motif (QxVxG)

110 were removed from posterior analysis.

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

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

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

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

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

116    MODEL program (http://swissmodel.expasy.org/interactive). The known crystal structure of rice

117    oryzacystatin I (OC-I) and SiCYS (Hu et al. 2015; Yuan et al. 2016) were used to construct the

118    homology-based models. Structure analysis was conducted by the RasMol 2.7 program.

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

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

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

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

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

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

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

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

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

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

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

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

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

132    less than 100 kb single region or separated by less than five genes. For microsynteny analysis,

133    the CDS sequence of every cystatin from *Arabidopsis*, barley, rice, and Sorghum was used as the

134    query to search against all other cystatins using NCBI_blast software with *e*-value $\leq$ 1e$^{-10}$. The

135    Circos software was used to display the results of collinear gene pairs (Krzywinski et al. 2009).

136    **Calculation of Ka and Ks**

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

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

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

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

141    **Promoter analysis of *SbCys* genes**

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

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

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

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

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

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

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

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

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

151    homologous relationship. The interaction network of SbCys proteins was drawn by

152    Cytoscape_v3.6.0.

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

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

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

156    (Wang et al. 2018). Root, shoot, and seedling were collected at 14 days after germination.

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

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

159    10 mm, and 1-2 cm. Three biological replicates were performed for each plant tissue. RNA was

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

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

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

163  infected with *Bipolaris sorghicola* (Yazawa et al. 2013). RNA samples were collected at 0, 12,

164  and 24 hours post-inoculation with one biological replicate. RNA-seq was run using Illumina

165  technology to give 100-base-pair single-end reads on a HiSeq2000 system. The second study

166  measured changes in the transcriptome of Sorghum leaves infested by sugarcane aphid (Tetreault

167  et al. 2019). The RNA-seq data were downloaded from the NCBI SRA database. In this study,

168  two treatments (infested and control) were arranged and two Sorghum genotypes (resistant

169  cultivar RTx2783 and susceptible cultivar BCK60) were used. Leaf samples were collected from

170  treated and control plants at 5, 10, and 15 days post sugarcane aphid infestation. Three biological

171  replicates were performed for all treatment and time combinations. RNA was sequenced using

172  the Illumina Hiseq 2500 platform to generate 100 bp single end reads. The accession numbers

173  and sample information were listed in Table S1. The differential expression of *SbCys* genes were

174  investigated by Hisat2 (http:/kim-lab.org/), Htseq (http://www.htseq.org/), and DESeq2 (R

175  package) based on the RNA-seq data (Wen, 2017). The $p \leq 0.05$ and $|logFC| \geq 1.5$ were set as the

176  cut-off criterion.

177  **Plant materials and treatments**

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

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

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

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

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

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

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

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

186   2017). The plants were transferred quickly to the nutrient solution containing 0.1 mM ABA

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

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

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

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

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

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

193   containers were used for three biological replicates respectively.

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

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

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

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

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

199   transcription was performed according to the kit instructions (Promega, Madison, USA). Primer

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

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

202   0.4 μl of each primer (forward and reverse), 2 μl 10-fold diluted cDNA, 7.2 μl of nuclease-free

203   water, and 10 μl of GoTaq® qPCR Master Mix (Perfect Real Time; Promega). PCR reaction

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

205   72 ˚C for 20s. The reactions were conducted using the CFX96 Real-Time PCR Detection System

206   (Bio-Rad  Laboratories,  Inc.).  Three  independent  biological  replicates  and  two  technical

207 replicates of each sample were performed. Gene-specific amplification of both reference and

208 *cystatin* genes were standardized by the presence of a single, dominant peak in the qRT-PCR

209 dissociation curve analyses. All data were analyzed by CFX Manager Software (Bio-Rad

210 Laboratories, Inc.). The efficiency range of the qRT-PCR amplifications for all of the genes

211 tested was between 91% and 100%. The average target (*SbCys*) cT (threshold cycle) values were

212 normalized to reference (*β-actin*) cT values. The fold change between treated sample and control

213 was calculated using the slightly modified $2^{-(\Delta\Delta Ct)}$ method as described by Kebrom et al. (2010).

214 A probability of $p \leq 0.05$ was considered to be significant.

215

216 **RESULTS**

217 **Identification and analysis of *SbCys* genes**

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

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

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

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

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

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

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

225 numbers, protein sequences, and annotations assigned to GO terms of the identified SbCys

226 proteins were listed in Table S3. The average length of these SbCys proteins was 148 amino acid

227 residues and the length mainly centered on the range of 105 to 240 amino acid residues.

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

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

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

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

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

233    The analysis of exon-intron structure can provide useful 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 affect the gene structure (Altenhoff et al. 2012).

240    **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. These motifs formed a fundamental structural

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

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

250    role in the inhibitory capacity of cystatins towards their target cysteine proteases (Meriem et al.

251    2010). Motif 2 contained a particular consensus sequence ([LVI][GA][RQG][WF]AV) that

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

253    typical motifs for SbCys proteins, motif 3 (V[WY][EVG]KPW) and motif 4 ([RK]xLxxF),

254    which were firstly described in tobacco (Zhao et al. 2014), were also detected in most SbCys

255    proteins, indicating their conserved and common role in both dicots and monocots. Motif 5

256    existed only in 3 SbCys family members (SbCys5, SbCys8, and SbCys15). Details of the 5

257    conserved motifs were shown in Fig. 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, might due to an extra α-helix that existed

264    in the C-terminal extension of SbCys14. Two important motifs (the conserve QxVxG motif and

265    W residue) of Sorghum cystatins involved in the interaction with the target cysteine enzymes

266    were also shown in Fig. 4. The predicted structure of SbCys13 showed some distortions in the

267    region of the β2 sheet, probably due to the insertion of a methionine in the first position of the

268    conserved QxVxG motif.

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

270    The cystatin gene family is highly conserved in both monocots and dicotyledons (Martinez and

271    Diaz, 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 related to rice and barley than *Arabidopsis*. Interestingly, a pair of *SbCys* genes (*SbCys2-*

283    *1* and *SbCys2-2*) was involved in the tandem duplication event in Sorghum (Fig. S2). Analysis of

284    duplicated *SbCys* genes showed that the *Ka*/*Ks* ratio far less than 1, varying from 0.0976 to

285    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 the biotic and abiotic stresses. All

295    of *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, HSE element was located in 10 *SbCys* genes, and TC-rich repeats

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

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

300  storage protein gene expression (Strömvik and Fauteux, 2009). The high diversity of the *cis*-

301  acting elements suggested that these *SbCys* genes might have a wide range of functional roles

302  and could be involved in multiple stress responses and growth and development progress (Zhang

303  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, seedling, pollen, endosperm, embryo, inflorescence (1-5mm, 1-

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

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

323   expression patterns of *SbCys* genes were significantly different between reproductive tissues and

324   vegetative tissues, such as *SbCys2-1*, *SbCys3*, *SbCys4*, *SbCys5*, *SbCys7*, *SbCys9*, *SbCys12*, and

325   *SbCys17*, which showed relatively higher expression levels in reproductive tissues including

326   pollen, endosperm, embryo, and pericarp than in vegetative tissues, while the expression of

327   *SbCys7* and *SbCys15* were higher in vegetative tissues than in reproductive tissues. It was worth

328   noting that the majority of *SbCys* genes had lower expression levels during inflorescence

329   development excepting *SbCys17* which displayed a higher expression pattern.

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

331   To gain insight into the potential roles of *SbCys* genes in response to *Bipolaris sorghicola*

332   infection and sugarcane aphid infestation, the relative expression patterns of these genes were

333   investigated by using the public transcription data from NCBI SRA database (DRP000986 and

334   SRP162227, respectively). As shown in Fig. 9 and 10, the expression patterns of *SbCys* genes

335   were different under the two biotic stresses. In response to *Bipolaris sorghicola* infection, seven

336   *SbCys* genes were induced and only 2 genes (*SbCys12* and *SbCys13*) were suppressed in the

337   infected Sorghum leaves compared with control (Fig. 9a). However, under aphid infestation, four

338   *SbCys* genes (*SbCys4*, *SbCys10*, *SbCys11*, and *SbCys14*) were up-regulated and 3 genes (*SbCys1*,

339   *SbCys3*, and *SbCys17*) were down-regulated relative to control in the susceptible Sorghum line

340   (BCK60). In the resistant Sorghum line (RTx2783), only two *SbCys* genes (*SbCys4* and *SbCys11*)

341   were induced, and the rest were barely expressed in Sorghum leaves with aphid infestation (Fig.

342   9b and 10). These results might suggest that *SbCys* genes played different roles in responding to

343   pathogen infection and aphid infestation.

344   **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-down expression trends, but *SbCys5* displayed a down-up expression

355    pattern (Fig. 11b). After exogenous ABA treatment, the expressions of 4 *SbCys* genes (*SbCys2-2*,

356    *SbCys3*, *SbCys4*, and *SbCys7*) were significantly up-regulated at three time points, but 9 genes

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

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

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

360    stresses except *SbCys4* that was significantly induced under dehydration, salt, and ABA stresses,

361    suggesting that SbCys4 might play an important role in response to different stress responses.

362

363    **DISCUSSION**

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

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

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

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

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

369  2015; Dutt et al. 2016; Yuan et al. 2016). However, little is known about the cystatin family in

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

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

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

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

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

375  to environment.

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

377  of them were distributed on chromosome 1 (Fig. 1). The uneven distribution of *cystatin* genes in

378  chromosomes was also found in the *B*. *distachyon* genome and the *Oryza sativa* genome

379  (Subburaj et al. 2017; Wang et al. 2015). This phenomenon might be due to the tandem

380  duplication events of *cystatin* genes on the chromosomes (Li et al. 2017). Several tandem

381  duplication events occurred at chromosomes 1 of the *B*. *distachyon* genome (Subburaj et al.

382  2017). Two tandem duplication events (*OsCys4*/*OsCys5* and *OsCys6*/*OsCys7*) were found

383  among *OsCys* genes, and existed in chromosomes 1 and 3 (Wang et al. 2015). One tandem

384  duplication event (*SbCys2-1*/*SbCys2-2*) occurred at chromosome 1 of the Sorghum genome (Fig.

385  S2). Eighteen *SbCys* genes were divided into three groups based on phylogenetic analysis (Fig.

386  5). Some conserved motifs among SbCys proteins had been identified by the alignment of the

387  amino acid sequences (Fig. 3). However, the conservation was accompanied with the differences

388  in some important amino acids indicated that SbCys family members might undergo a complex

389  evolutionary history. The variation of crucial amino acids of cystatins might have a significant

390  influence on their respective functions (Tremblay et al. 2019). For example, the QxVxG motif

391   could directly enter and interact with the active site of targeted enzymes. The motif was

392   conserved in all SbCys proteins with the exceptions of 5 cystatins (SbCys1, SbCys6, SbCys8,

393   SbCys9, and SbCys13) that were partially modified by the insertion or variation in important

394   residues (Fig. 3a). Furthermore, three SbCys proteins (SbCys8, SbCys9, and SbCys13) showed

395   significant variations with other Sorghum cystatins in their predicted three-dimensional

396   structures (Fig. 4). The variations in vital amino acid residues might result in the change in

397   cystatin inhibitory action (Tremblay et al. 2019). In addition, two novel motifs, motif 3

398   (V[WY][EVG]KPW) and motif 4 ([RK]xLxxF), firstly described in tobacco (Zhao et al. 2014),

399   were also identified in the C-terminalin of many SbCys proteins. The contribution of the two

400   new motifs to cystatin inhibitory action needs to be further studied.

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

402   endogenous and exogenous cysteine proteases activities during seed development (Tremblay et

403   al. 2019). In the present study, as revealed by RNA-seq data analysis (Fig. 8 and S3), the

404   expression levels of several *SbCys* family genes were higher in reproductive tissues than in

405   vegetative tissues, which were consistent with the reports that most cystatins were specifically

406   expressed in developing seeds and played a role in seed development (Dutt et al. 2010; Zhao et al.

407   2014). Moreover, promoter analysis showed that the highly expressed *SbCys* genes in

408   reproductive tissues possessed endosperm expression-related *cis*-elements (Skn-1 and

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

410   several SbCys proteins could interact with many functional proteins (e.g., growth and

411   development related proteins, vacuolar-processing enzyme, and lipid-transfer proteins) (Fig. 7),

412   implying these cystatins were involved in regulating the gene expression of cereal grain storage

413   proteins (Diaz-Mendoza et al. 2016).

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

415    proteins against pest infestation and pathogen infection (Meriem et al. 2010). At present, some

416    cystatins with insecticidal activity have been isolated from many plants, such as barley, tomato,

417    and potato (Rasoolizadeh et al. 2017; Siddiqui et al. 2017; Velasco-Arroyo et al. 2018; Goulet et

418    al. 2020). Several cystatins having antifungal activities were also isolated from taro, cacao, and

419    wheat (Christova et al. 2018; Pirovani et al. 2010; Chen et al. 2014). Although studies on

420    insecticidal and antifungal activity of plant cystatins have been well established *in vitro*, the

421    knowledge about their roles in plants in response to biotic stresses is limited. To explore the

422    properties of *SbCys* genes responding to pest infestation and pathogen infection, we conducted

423    the analysis on the expression patterns of *SbCys* genes. The results showed that the expressions

424    of most *SbCys* genes were induced during *Bipolaris sorghicola* infection, suggesting these

425    cystatins played functions in inhibiting exogenous cysteine proteases secreted by pathogens to

426    infect plant tissues (Fig. 9a). Interestingly, for sugarcane arthropods infestation, only two genes

427    (*SbCys4* and *SbCys11*) were up-regulated significantly in susceptible and resistant Sorghum lines

428    (Fig. 9b and 10). These differential expression patterns between *SbCys* genes might suggest that

429    some of them had evolved to inhibit specific cysteine proteinases. The exact roles of these *SbCys*

430    genes in insecticidal and antifungal activity *in vivo* are worthy to be explored in further study.

431    *Cystatin* genes are also involved in various abiotic stress responses in plants. In *Arabidopsis*, the

432    expression levels of *AtCYS1* and *AtCYS2* were enhanced by high temperature and wounding

433    stresses (Hwang et al. 2010). *AtCYSa* and *AtCYSb* were also induced by different abiotic stresses,

434    e.g., salt, drought, oxidation, and cold stresses (Zhang et al. 2008). Velasco-Arroyo et al. (2018)

435    reported that the silence of barley *HvCPI-2* and *HvCPI-4* specifically modified leaf responses to

436    drought stress. Wang et al. (2015) observed the significant change in the expression levels of

437  several rice *OsCYS* genes under cold, drought, salt, and hormone treatments. In the present study,

438  most *SbCys* genes were found to have positive or negative responses to dehydration, salt, and

439  ABA stresses. Moreover, the interaction results showed that most cystatins could interact with

440  stresses-related proteins, implying that the cystatins played critical roles in response to diverse

441  stress conditions. Notably, the expression of *SbCys4* was significantly up-regulated under three

442  stress conditions (Fig. 11), suggesting a specific role of SbCys4 in responding to various stress

443  conditions. Promoter analysis indicated that stress-related *cis*-elements were widespread in the

444  promoter region of these cystatin genes (Table S5), and *SbCys4* possessed plenty of stress-related

445  *cis*-elements, including G-box, ABRE, HSE, MBS, and TC-rich repeats. These results provide an

446  effective reference for the functional verification of the *SbCys* family genes under abiotic

447  stresses.

448

449  **CONCLUSIONS**

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

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

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

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

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

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

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

457  Furthermore, the variation of amino acids in Sorghum cystatin critical active sites suggested that

458  they might undergo a complex evolutionary process and possess structural and functional

459  divergence. The expression profiles of *SbCys* genes in different tissues indicated that most *SbCys*

460  genes were involved in plant growth and development. Changes in the expression of *SbCys*

461  genes under biotic and abiotic stresses indicated that many *SbCys* genes played important roles in

462  response to unfavorable growth conditions. It was worth noting that the expression of *SbCys4*

463  was significantly enhanced under biotic and abiotic stresses, suggesting its unique role in

464  mediating the response of Sorghum to adverse environmental conditions.

465

466  **REFERENCES**

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

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

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

470      10.1371/journal.pcbi.1002514.

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

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

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

474      10.1046/j.1432-1033.2003.03630.x.

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

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

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

478      10.1111/pce.13178.

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

501    **Goulet MC, Sainsbury F, Michaud D. 2020.** Cystatin activity-based protease profiling to select

502         protease inhibitors useful in plant protection. *Methods in Molecular Biology* **2139:**353-366

503         DOI 10.1007/978-1-0716-0528-8.

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

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

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

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

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

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

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

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

512            **2010.** Distinct expression patterns of two *Arabidopsis* phytocystatin genes, AtCYS1 and

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

514            10.1007/s00299-010-0876-y.

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

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

517            **33(1):**48-58 DOI 10.4161/psb.5.3.11186.

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

519            **2010.** Deleterious effects of plant cystatins against the banana weevil *Cosmopolites*

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

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

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

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

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

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

526            **2009.** Circos: an information aesthetic for comparative genomics. *Genome research*

527            **19(9):**1639-1645 DOI 10.1101/gr.092759.109.

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

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

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

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

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

533        10.3390/genes8100290.

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

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

536        10.2174/0929866521666140418101711.

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

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

539        10.1186/s12864-015-1554-9.

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

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

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

543        **151(3):**1531-1545 DOI 10.1104/pp.109.146019.

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

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

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

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

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

549        **8(1):**198-210 DOI 10.1186/1471-2148-8-198.

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

551        **2016.** Phytocystatins: defense proteins against phytophagous insects and acari.

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

553 10.3390/ijms17101747.

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

555 **92(11):**1657-1666 DOI 10.1016/j.biochi.2010.06.006.

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

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

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

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

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

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

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

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

564 **232(6):**1485-1497 DOI 10.2307/23391912.

565 **Rasoolizadeh A, Goulet MC, Guay JF, Cloutier C, Michaud D. 2017.** Population-associated

566 heterogeneity of the digestive Cys protease complement in Colorado potato beetle,

567 *Leptinotarsa decemlineata. Journal of Insect Physiology* **106:**125-133 DOI

568 10.1016/j.jinsphys.2017.03.001.

569 **Siddiqui AA, Khaki PS, Bano B. 2017.** Interaction of almond cystatin with pesticides:

570 Structural and functional analysis. *Journal Molecular Recognitnition* **30(3):**e2586 DOI I

571 10.1002/jmr.2586.

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

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

574 **40(8):**577-586 DOI 10.14348/molcells.2017.0075.

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

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

577         10.1186/1471-2229-9-126.

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

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

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

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

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

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

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

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

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

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

588         10.1016/j.plaphy.2017.03.025.

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

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

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

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

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

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

595         10.3389/fpls.2019.00145.

596    **Tremblay J, Goulet MC, Michaud D. 2019.** Recombinant cystatins in plants. *Biochimie*

597         **166:**184-193 DOI 10.1016/j.biochi.2019.06.006.

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

599          Multifunctional amaranth cystatin inhibits endogenous and digestive insect cysteine

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

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

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

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

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

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

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

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

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

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

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

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

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

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

614          **34(9):**1579-1592 DOI 10.1007/s00299-015-1810-0.

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

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

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

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

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

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

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

622        10.1093/dnares/dsu039.

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

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

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

626        10.1016/j.plantsci.2017.07.002.

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

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

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

630        10.1371/journal.pone.0062460.

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

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

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

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

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

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

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

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

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

640        DOI 10.1093/jxb/eru274.

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

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

643        *Bioinformatics* **4(4):**259-263 DOI 10.1016/S1672-0229(10)60008-3.
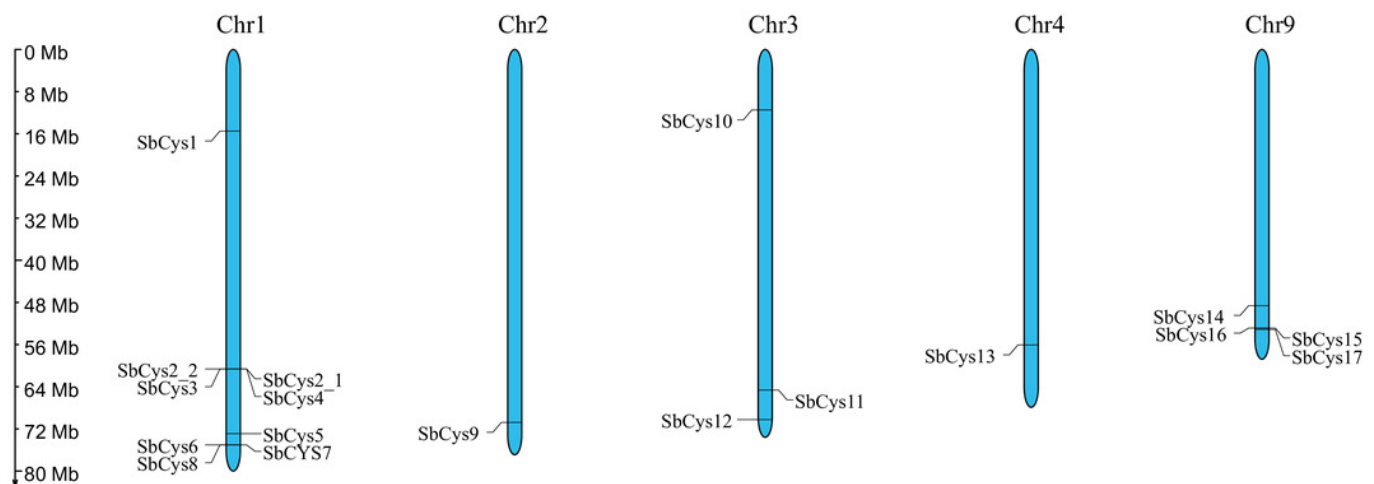
644

645

646

647

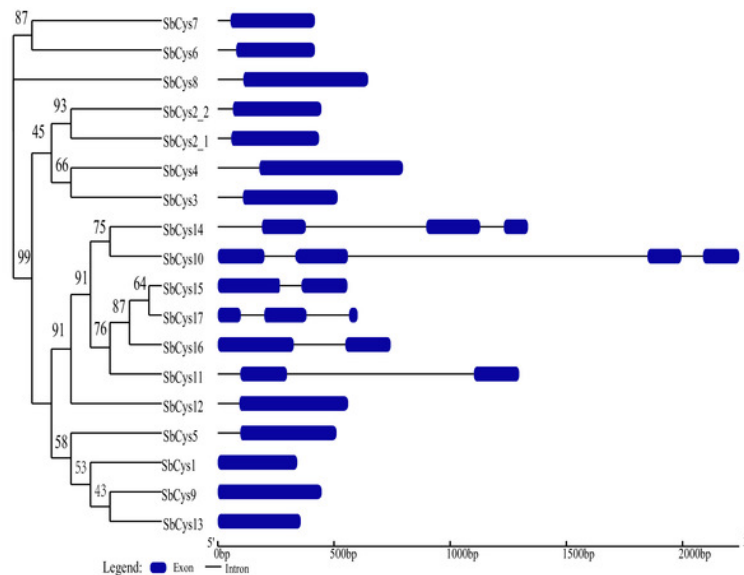# Figure 1

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

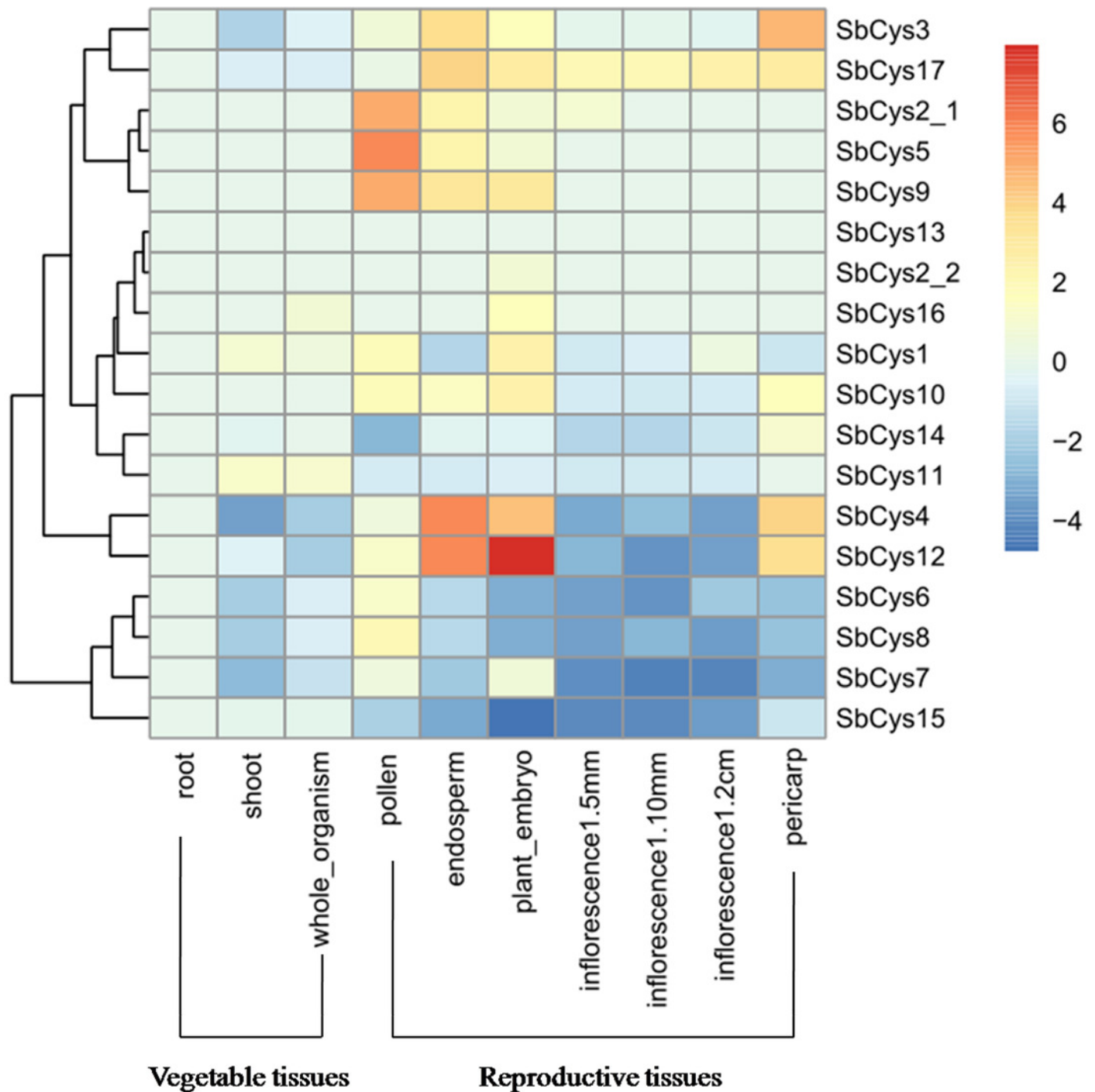The interaction networks of SbCys proteins according to the orthologs in *Arabidopsis*.

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.
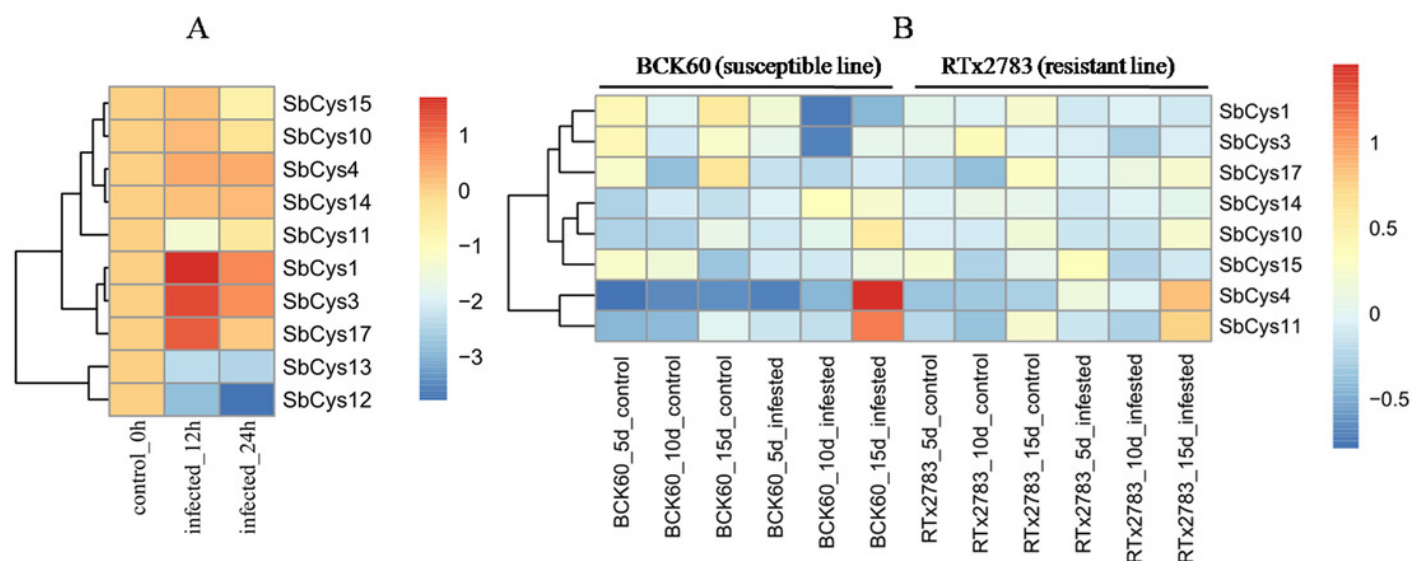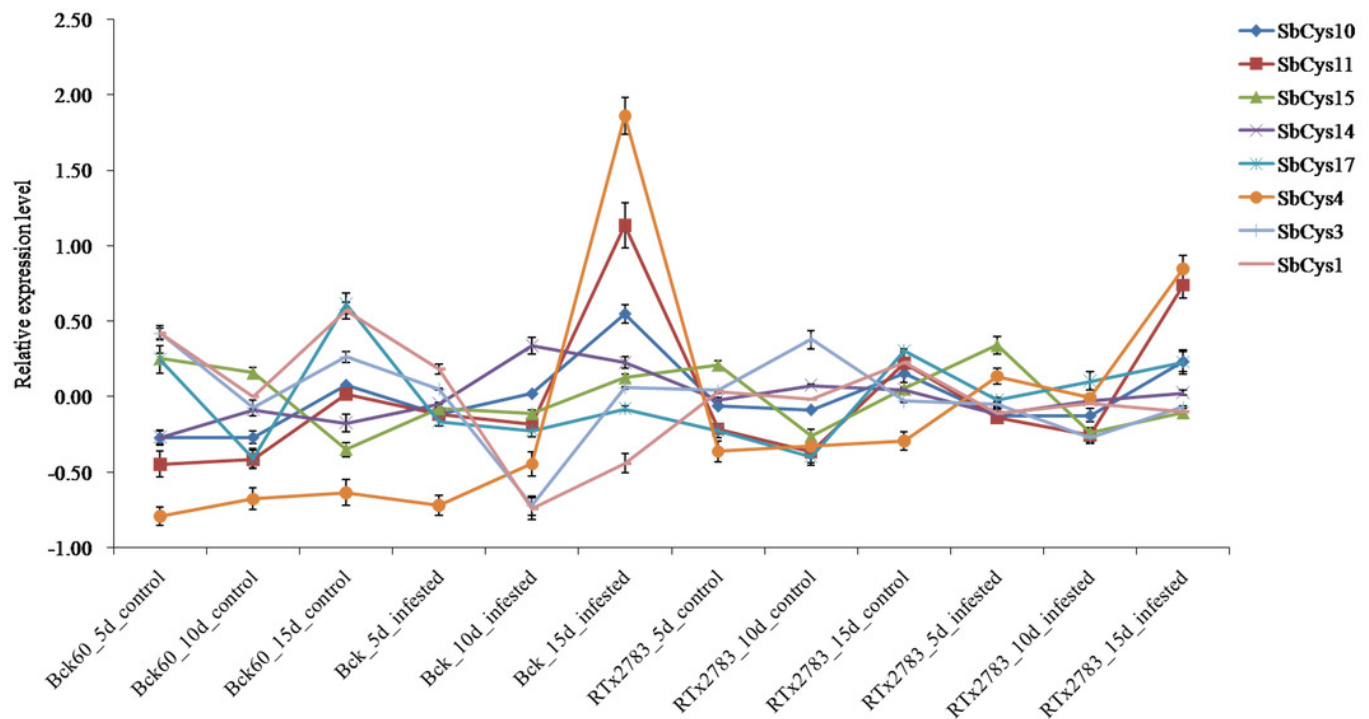
# Figure 9

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 $p \leq 0.01$.