

**A peer-reviewed version of this preprint was published in PeerJ on 7 September 2016.**

[View the peer-reviewed version](https://doi.org/10.7717/peerj.2437) (peerj.com/articles/2437), which is the preferred citable publication unless you specifically need to cite this preprint.

Bi C, Xu Y, Ye Q, Yin T, Ye N. 2016. Genome-wide identification and characterization of WRKY gene family in *Salix suchowensis*. PeerJ 4:e2437 <https://doi.org/10.7717/peerj.2437>

# Genome-wide identification and characterization of WRKY gene family in *Salix suchowensis*

Changwei Bi<sup>1</sup>, Yiqing Xu<sup>1</sup>, Qiaolin Ye<sup>1</sup>, Tongming Yin<sup>2</sup>, Ning Ye<sup>Corresp. 1</sup>

<sup>1</sup> College of Information Science and Technology, Nanjing Forestry University, Nanjing, Jiangsu, China

<sup>2</sup> College of Forest Resources and Environment, Nanjing Forestry University, Nanjing, Jiangsu, China

Corresponding Author: Ning Ye  
Email address: yening@njfu.edu.cn

WRKY proteins are the plant-specific zinc finger transcription factors. They can specifically interact with the W-box ([C/T]TGAC[T/C]), which can be found in the promoter region of a large number of plant target genes, to regulate the expressions of downstream target genes. They also participate in diverse physiological and growing processes in plants. Prior to the present studies, plentiful WRKY genes have been identified and characterized in herbaceous species, but there is no large-scale study of WRKY genes in willow. With the whole genome sequencing in *Salix suchowensis*, we have the opportunity to conduct the genome-wide research for willow WRKY gene family. In this study, we identified 85 WRKY genes in the willow genome and renamed them from SsWRKY1 to SsWRKY85 on the basis of their specific distributions on chromosomes. Due to their diverse structural features, the 85 willow WRKY genes could be further classified into three main groups (group I - III), with five subgroups (IIa - IIe) in group II. With the multiple sequence alignment and the manual search, we found three variations of the WRKYGQK heptapeptide: WRKYGRK, WKKYGQK and WRKYGKK, and four variations of the normal zinc finger motif, which might execute some new biological functions. In addition, the SsWRKY genes from the same subgroup share the similar exon-intron structures and conserved motif domains. Further studies of SsWRKY genes revealed that segmental duplication events played the prominent roles in the expansion of SsWRKY genes. Distinct expression profiles of SsWRKY genes with RNA sequencing data revealed that diverse expression patterns among five tissues, including tender roots, young leaves, vegetative buds, non-lignified stems and barks. With the analyses of WRKY gene family in willow, it is not only beneficial to complete the functional and annotation information of WRKY genes family in woody plants, but also provide important references to investigate the expansion and evolution of this gene family in flowering plants.

# 1 **Genome-wide identification and characterization of** 2 **WRKY gene family in *Salix suchowensis***

3 Changwei Bi<sup>1</sup>, Yiqing Xu<sup>1</sup>, Qiaolin Ye<sup>1</sup>, Tongming Yin<sup>2</sup>, Ning Ye<sup>1\*</sup>

4 1. College of Information Science and Technology, Nanjing Forestry University, Nanjing,  
5 Jiangsu, China

6 2. College of Forest Resources and Environment, Nanjing Forestry University, Nanjing,  
7 Jiangsu, China

8 \* Corresponding author, yening@njfu.edu.cn

9 Changwei Bi, bichwei@163.com

10 Ning Ye, yening@njfu.edu.cn

## 11 **Abstract**

12 WRKY proteins are the plant-specific zinc finger transcription factors. They can specifically  
13 interact with the W-box ([C/T]TGAC[T/C]), which can be found in the promoter region of a  
14 large number of plant target genes, to regulate the expressions of downstream target genes.  
15 They also participate in diverse physiological and growing processes in plants. Prior to the  
16 present studies, plentiful WRKY genes have been identified and characterized in herbaceous  
17 species, but there is no large-scale study of WRKY genes in willow. With the whole genome  
18 sequencing in *Salix suchowensis*, we have the opportunity to conduct the genome-wide  
19 research for willow WRKY gene family. In this study, we identified 85 WRKY genes in the  
20 willow genome and renamed them from SsWRKY1 to SsWRKY85 on the basis of their  
21 specific distributions on chromosomes. Due to their diverse structural features, the 85 willow  
22 WRKY genes could be further classified into three main groups (group I - III), with five  
23 subgroups (IIa - IIe) in group II. With the multiple sequence alignment and the manual search,  
24 we found three variations of the WRKYGQK heptapeptide: WRKYGRK, WKKYGQK and

WRKYGKK, and four variations of the normal zinc finger motif, which might execute some new biological functions. In addition, the SsWRKY genes from the same subgroup share the similar exon–intron structures and conserved motif domains. Further studies of SsWRKY genes revealed that segmental duplication events played the prominent roles in the expansion of SsWRKY genes. Distinct expression profiles of SsWRKY genes with RNA sequencing data revealed that diverse expression patterns among five tissues, including tender roots, young leaves, vegetative buds, non-lignified stems and barks. With the analyses of WRKY gene family in willow, it is not only beneficial to complete the functional and annotation information of WRKY genes family in woody plants, but also provide important references to investigate the expansion and evolution of this gene family in flowering plants.

**Keywords:** WRKY, Phylogenetic analysis, Evolution, Duplication, Expression, Willow

## Introduction

Plants form a series of adjustment mechanisms to adapt diverse environment stress in their long evolutionary processes. Among the numerous adjustment mechanisms, transcription factors play important roles [1]. In plants, WRKY proteins constitute a large family of transcription factors, involving in various physiological and developmental processes [2, 3]. Since the first WRKY gene was cloned and characterized from sweet potato [4], many corresponding studies have been conducted rapidly, such as *Arabidopsis thaliana*, desert legume (*Retama raetam*), cotton (*Gossypium arboreum*), rice (*Oryza sativa*), *Pinus monticola*, barley (*Hordeum vulgare*), sunflower, cucumber (*Cucumis sativus*), poplar (*Populus trichocarpa*), tomato (*Solanum lycopersicum*) and grapevine (*Vitis vinifera*) [2, 5-14].

The existence of either one or two highly conserved WRKY domains is the most vital structural characteristic of WRKY gene. WRKY gene consists of about 60 amino acid residues with a conserved WRKYGQK heptapeptide at its N-termini, and a zinc finger motif (C-X<sub>4-5</sub>-C-X<sub>22-23</sub>-H-X<sub>1</sub>-H or C-X<sub>7</sub>-C-X<sub>23</sub>-H-X<sub>1</sub>-C) at the C-terminal region. Previous functional studies indicated that WRKY genes could specifically interact with the W-box, the promoter region of plant target genes, to adjust the expressions of downstream target genes

[15]. What's more, SURE (sugar responsive elements), another prominent cis-element that can promote transcription processes, was also found to bind to the WRKY transcription factors under a convincing research [16]. The proper DNA-binding ability of WRKY genes could be influenced by the variation of the conserved WRKYGQK heptapeptide [17, 18].

The WRKY proteins can be classified into three main groups (I, II and III) on the basis of the number of their WRKY domains and the pattern of the zinc finger motif. Proteins from group I contain two WRKY domains followed by a C<sub>2</sub>H<sub>2</sub> zinc finger motif, while the other WRKY proteins from group II and III only contain one WRKY domain followed by a C<sub>2</sub>H<sub>2</sub> or C<sub>2</sub>HC correspondingly [19]. Group II can be further divided into five subgroups from IIa to IIe based on additional amino acid motifs present outside the WRKY domain. Apart from the conserved WRKY domains and the zinc finger motif, there are also some WRKY proteins appearing to have basic nuclear localization signal, LZs (leucine zipper) [20], serine-threonine-rich region, glutamine-rich region and proline-rich region [21]. Throughout the studies of WRKY gene family in many higher plants [3, 10, 13], WRKY genes have been identified to be involved in various regulatory processes mediated by different biotic and abiotic stresses [22]. In plant defense against various biotic stresses, such as bacterial, fungal and viral pathogens, it has been well documented that the WRKY genes play vital roles [14, 23, 24]. They are also involved in abiotic stress-induced gene expression. In *Arabidopsis*, with the either heat or salt treatments, the expressions of AtWRKY25 and AtWRKY33 are transformed apparently [25]. Furthermore, the expression of TcWRKY53 that belonged to alpine penny grass (*Thlaspi caerulescens*) is affected by salt, cold, and polyethylene glycol treatments [3]. In rice, a total of 54 OsWRKY genes showed noticeable differences in their transcript abundance under the abiotic stress such as cold, drought, and salinity [22]. There is also accumulating evidence that WRKY genes are involved in regulating developmental processes, such as embryo morphogenesis [26], senescence [27], trichome initiation [28], and some signal transduction processes mediated by plant hormones including gibberellic acid [29], abscisic acid [30], or salicylic acid [31].

The number of WRKY genes in different species varies tremendously. For instance, there are 72 members in *Arabidopsis thaliana*, at least 45 in barley, 57 in cucumber, 58 in physic nut (*Jatropha curcas*), 59 in grapevine, 104 in poplar, 105 in foxtail millet (*Setaria italica*), 112 in *Gossypium raimondii* and more than 109 in rice [2, 6, 7, 9, 11, 13, 32-34]. Zhang et al. also identified the most basal WRKY genes in the lineage of non-plant eukaryotes and green alga [35]. The study in bryophyte (*Physcomitrella patens*) found at least 12 WRKY genes [21], and the study in gymnosperm (*Cycas revolute*) identified at least 21 WRKY genes [36]. Interestingly, the WRKY genes in eukaryotic unicellular chlamydomonas, protocist (*Giardia lamblia*), bryophyte (*Physcomitrella patens*) and fern (*Ceratopteris richardii*) all belonged to group I [2, 37]. The WRKY genes in *Cycas revolute* were divided into two groups, 15 WRKY genes therein belonged to group I and the other 6 WRKY genes belonged to group II. Further study suggested that the core WRKY domains of group II and III were similar to the C-terminal domain of group I, and the group II WRKY genes might emerge from the breakage of the C-terminal domain in group I and the group III probably evolve from group III [21]. Above of all indicated that the group I WRKY genes might be the oldest type, which evolved from the origin of eucaryon, and group II and III might generate after the origin of bryophyte [35, 38]. In the evolution of WRKY genes, gene duplication events played prominent roles. As we all know, gene duplication events can lead to the generation of new genes. Take this an example, there are approximately 80% of OsWRKY (rice) genes located in duplicated regions [13], as well as 83% of PtWRKY (poplar) genes [7]. However, no gene duplication events have occurred in cucumber [9].

Willow, an important broad-leaf plant, grows quickly and reproduces simply. It can survive under a variety of different ecological environment and grow well. With its broad leaf, willow becomes a prominent part of the protection forest, soil and water conservation forest specie. Therefore, willow has higher ecological and economic value. With these various factors and the draft of the *Salix suchowensis* genome sequence was finished recently [39], we had the opportunity to analyze the willow WRKY gene family. In this study, we identified 85 members of the WRKY genes in the willow genome. Subsequently, the distribution of

1 WRKY genes on chromosomes, phylogenetic analysis, classification of WRKY genes,  
2 exon-intron organization, conserved motif analysis, and expression analyses were also  
3 conducted, which provide a solid foundation for further studies of SsWRKY gene family  
4 function and evolution.

## 5 **Materials and methods**

### 6 **Datasets and sequence retrieval**

7 The sequence of a shrub willow *Salix suchowensis* (*S. suchowensis*), which flowers within  
8 two years, was conducted with a combined approach using Roche/454 and  
9 Illumina/HiSeq-2000 sequencing technologies [39]. The latest v5.2 *S. suchowensis* genome  
10 annotation information (version5\_2.gff3) and protein sequences (Willow.gene.pep) were  
11 downloaded from our laboratory website ([http://bio.njfu.edu.cn/ss\\_wrky/](http://bio.njfu.edu.cn/ss_wrky/)). Sequences of 72  
12 *Arabidopsis* WRKY proteins were obtained from TAIR (release 10,  
13 <http://www.arabidopsis.org/>) [2], and 104 poplar WRKY proteins were obtained from the  
14 Supplementary material 3 of poplar [7].

### 15 **Identification and distribution of WRKY genes in willow**

16 The procedure performed to identify putative WRKY proteins in willow was similar to the  
17 method described in other species [6, 7, 13]. The Hidden Markov Model (HMM) profile for  
18 the WRKY transcription factor was downloaded from the Pfam database  
19 (<http://pfam.sanger.ac.uk/>) with the keyword 'PF03106' [40]. The HMM profile was applied  
20 as a query to search against the all willow protein sequences (Willow.gene.pep) using  
21 BLASTP program ( $E\text{-value} = 1e^{-3}$ ) [41]. Another procedure was performed to validate the  
22 putative accuracy. An alignment of WRKY seed sequences in Stockholm format from Pfam  
23 database was used by HMMER program (hmmbuild) to build a HMM model, and then the  
24 model was used to search the willow protein sequences by another HMMER program  
25 (hmmsearch) with default parameters [42]. Finally, we employed the SMART program

1 (<http://smart.embl-heidelberg.de/>) to confirm the candidates from the two procedures  
2 correlated with the WRKY structure features [43].

3 Additionally, we calculated the length, MW (molecular weight), PI (isoelectric point) of  
4 these putative WRKY proteins by ExPasy site ([http://au.expasy.org/tools/pi\\_tool.html](http://au.expasy.org/tools/pi_tool.html)). Every  
5 WRKY genes were mapped onto chromosomes assembled ourselves  
6 ([http://bio.njfu.edu.cn/ss\\_wrky/version5\\_2.fa](http://bio.njfu.edu.cn/ss_wrky/version5_2.fa)) with an in-house Perl script  
7 ([http://bio.njfu.edu.cn/willow\\_chromosome/BuildGff3\\_Chrom.pl](http://bio.njfu.edu.cn/willow_chromosome/BuildGff3_Chrom.pl)), and then rename based on  
8 their orderly given chromosomal distribution. The distribution graph of every WRKY gene  
9 was drawn by MapInspect software (<http://mapinspect.software.informer.com/>).

## 10 **Sequence alignments, phylogenetic analysis and classification of** 11 **willow WRKY genes**

12 Using the online tool SMART, we obtained the conserved WRKY core domains of predicted  
13 SsWRKY genes, and then multiple sequence alignment based on these domains was  
14 performed using ClustalX (version 2.1) [44]. After alignment, we used Boxshade  
15 ([http://www.ch.embnet.org/software/BOX\\_form.html](http://www.ch.embnet.org/software/BOX_form.html)) to color the alignment result online. To  
16 gain better classification of these SsWRKY genes, a further multiple sequence alignment  
17 including 103 SsWRKY domains and 82 WRKY domains from *Arabidopsis* (AtWRKY) was  
18 performed using ClustalW [44], and a phylogenetic tree based on this alignment was built by  
19 MEGA 6.0 with the Neighbor-joining (NJ) method [45]. Bootstrap values have been  
20 calculated from 1000 iterations in the pairwise gap deletion mode, which is conducive to the  
21 topology of the NJ tree by divergent sequences. Based on the phylogenetic tree constructed by  
22 SsWRKY and AtWRKY domains, these SsWRKY genes were classified into different groups  
23 and subgroups. In order to get a better comparison of WRKY family in Salicaceae, a  
24 phylogenetic tree including all SsWRKY domains and 126 WRKY domains from poplar  
25 (PtWRKY) was constructed with the similar method to *Arabidopsis*. Additionally, a  
26 phylogenetic tree based on full-length SsWRKY genes was also constructed to get a better  
27 classification. The ortholog of each SsWRKY gene in *Arabidopsis* and poplar was based on



1 the phylogenetic trees of their respective WRKY domains, and the members of group I  
2 WRKY genes were considered as orthologs unless the same phylogenetic relationship can be  
3 detected between N-termini and C-termini in the tree. Another method, BLAST-based method  
4 (Bi-direction best hit) [46], was used to verify the putative orthologous genes (e-value =  
5 1e-20).

## 6 **Evolutionary analysis of WRKY III genes in willow**

7 The group of WRKY III genes, only found in flowering plants, was considered as the  
8 evolutionary youngest groups, and played crucial roles in process of plant growth and  
9 resistance [7, 13]. Previous study of Zhang et al. held the opinion that duplications and  
10 diversifications were plentiful in WRKY III genes, and they appeared to have confronted  
11 different selection challenges [35]. Phylogenetic analysis of WRKY III genes was performed  
12 using MEGA6.0 with 65 WRKY III genes from *Arabidopsis* (AtWRKY), *Populus*  
13 (PtWRKY), grape (VvWRKY), willow (SsWRKY) and rice (OsWRKY). A NJ tree was  
14 constructed with the same method described before. Additionally, we estimated the  
15 non-synonymous (Ka) and synonymous (Ks) substitution ratio of SsWRKY III genes to verify  
16 whether selection pressure participated in the expansion of SsWRKY III genes. Each pair of  
17 these WRKY III protein sequences was first aligned using ClustalW. The alignments  
18 generated by ClustalW and the corresponding cDNA sequences were submitted to the online  
19 program PAL2NAL (<http://www.bork.embl.de/pal2nal/>) [47], which automatically calculates  
20 Ks and Ka by the codeml program in PAML [48].

## 21 **Analysis of exon-intron structure, gene duplication events and** 22 **conserved motif distribution of willow WRKY genes**

23 The exon-intron structures of the willow WRKY genes were obtained based on the protein  
24 annotation files which we assembled ourselves  
25 ([http://bio.njfu.edu.cn/ss\\_wrky/version5\\_2.gff3](http://bio.njfu.edu.cn/ss_wrky/version5_2.gff3)), and the diagrams were obtained from the  
26 online website Gene Structure Display Server (GSDS: <http://gsds.cbi.pku.edu.cn/>) [49].

1 Gene duplication events were always considered as the vital sources of biological evolution.  
2 Blastp (e-value, 1e-20) was performed to identify the gene duplication events in SsWRKY  
3 genes with the following definition [7, 50]: (1) the coverage of the aligned sequence  $\geq 80\%$  of  
4 the longer gene; and (2) the similarity of the aligned regions  $\geq 70\%$ .

5 To better exhibit the structural features of SsWRKY proteins, the online tool MEME  
6 (Multiple Expectation Maximization for Motif Elicitation) was used to identify the conserved  
7 motifs in the encoded SsWRKY proteins [51]. The optimized parameters were employed as  
8 the following: any number of repetitions, maximum number of motifs = 20, and the optimum  
9 width of each motif was constrained to between 6 to 50 residues. The online program 2ZIP  
10 (<http://2zip.molgen.mpg.de/>) was used to verify the existence of the conserved Leu zipper  
11 motif [52], whereas some other important conserved motifs, HARF, LXXLL (X, any amino  
12 acid) and LXLXLX, were identified manually.

## 13 **Expression analyses of willow WRKY genes**

14 The sequenced *S. suchowensis* RNA-HiSeq reads from five tissues including tender roots,  
15 young leaves, vegetative buds, non-lignified stems and barks were separately mapped back  
16 onto the SsWRKY gene sequences using BWA (mismatch  $\leq 2$  bp, other parameters as  
17 default) [53], and the number of mapped reads for each WRKY gene was counted.  
18 Normalization of the mapped reads was done using RPKM (reads per kilo base per million  
19 reads) method [54]. The heat map for tissue-specific expression profiling was generated based  
20 on the  $\log_2$ RPKM values for each gene in all the tissue samples using R package [55].

## 21 **Results**

### 22 **Identification and characterization of 85 WRKY genes in willow** 23 **(*Salix suchowensis*)**

24 In this study, we obtained 92 putative WRKY genes by using HMMER to search the Hidden  
25 Markov Model profile of WRKY DNA-binding domain against willow protein sequences,

1 and validated the accuracy of the consequence by BlastP. After submitting the 92 putative  
2 WRKY genes to the online program SMART, seven genes without a complete WRKY  
3 domain were removed (willow\_GLEAN\_10004672, willow\_GLEAN\_10009126,  
4 willow\_GLEAN\_10011436, willow\_GLEAN\_10011470, willow\_GLEAN\_10018393,  
5 willow\_GLEAN\_10019671 and willow\_GLEAN\_10024347), and the other 85 WRKY genes  
6 were selected as possible members of the WRKY superfamily.

7 WRKY genes contain one or two WRKY domains, comprising a conserved WRKYGQK  
8 heptapeptide at the N-termini and a novel zinc finger motif (C-X<sub>4-7</sub>-C-X<sub>22-23</sub>H-X-H/C) at the  
9 C-termini [2]. The variations of WRKY core domain or zinc finger motif may lead to the  
10 binding specificities of WRKY genes, but this remains to be largely demonstrated [19, 56, 57].  
11 In order to identify the variations in WRKY core domains, a multiple sequence alignment of  
12 85 SsWRKY core domains was conducted, and the result was shown in Fig. 1. Among the  
13 selected 85 WRKY genes, 81 (95.3%) were identified to have highly conserved sequence  
14 WRKYGQK, whereas the other four WRKY genes (SsWRKY14, SsWRKY23, SsWRKY38  
15 and SsWRKY78) had a single mismatched amino acid in their core WRKY domains (Fig. 1).  
16 In SsWRKY14 and SsWRKY38, the WRKY domain has the sequence WRKYGKK, while  
17 SsWRKY23 contains a WKKYGQK sequence, and SsWRKY78 contains WRKYGRK  
18 sequence. Eulgem et al. previously described that the zinc finger motif (C-X<sub>4-5</sub>-X<sub>22-23</sub>-H-X<sub>1</sub>-H  
19 or C-X<sub>7</sub>-C-X<sub>23</sub>-H-X<sub>1</sub>-C) is another vital features of the WRKY family [2]. As illustrated in  
20 Fig. 1, four WRKY domains (SsWRKY76C, SsWRKY64, SsWRKY12 and SsWRKY28) do  
21 not contain any distinct zinc finger motif, but they were still reserved in the succeeding  
22 analyses, as performed in barley and poplar [7, 11]. Additionally, some zinc-finger-like motifs,  
23 including C-X<sub>4</sub>-C-X<sub>21</sub>-H-X<sub>1</sub>-H in SsWRKY23 and C-X<sub>5</sub>-C-X<sub>19</sub>-H-X<sub>1</sub>-H in SsWRKY73 and  
24 SsWRKY17, were identified in willow WRKY genes. Both the two zinc-finger-like motifs  
25 were also found in poplar (PtWRKY39, 57, 42 and 53).

26 Detailed characteristics of SsWRKY genes are list in Table 1, including the WRKY gene  
27 specific group numbers, chromosomal distribution, *Arabidopsis* and poplar orthologs. The  
28 molecular weight (MW), isoelectric point (PI) and the length of each WRKY protein

1 sequence are also shown in Table 1. According to the particularization (Table 1), the average  
2 length of these protein sequences is 407 residues, and the lengths ranged from 109 residues  
3 (SsWRKY23) to 1,593 residues (SsWRKY78). Additionally, the isoelectric point (PI) ranged  
4 from 5.03 (SsWRKY38, SsWRKY60) to 10.27 (SsWRKY28), and the molecular weight  
5 (MW) ranged from 12.9 (SsWRKY23) to 179.0 kDa (SsWRKY78).

## 6 **Locations and gene clusters of willow WRKY genes**

7 84 of the 85 putative SsWRKY genes could be mapped onto 19 willow chromosomes and  
8 then renamed from SsWRKY1 to SsWRKY84 based on their specific distributions on the  
9 chromosomes. Only one SsWRKY gene (willow\_GLEAN\_10002834), renamed as  
10 SsWRKY85, could not be conclusively mapped onto any chromosome. As shown in Fig. 2,  
11 Chromosome (Chr) 2 possessed the largest number of SsWRKY genes (11 genes), followed  
12 by Chr14 (10 genes). Eight SsWRKY genes were found on Chr6, six on Chr1 and Chr16, and  
13 five on Chr5. Additionally, four chromosomes (Chr4, Chr11, Chr17, Chr18) had four  
14 SsWRKY genes, as well as three SsWRKY genes were found on Chr8, Chr13 and Chr19.  
15 Chr10 and Chr15 had two SsWRKY genes, and only one SsWRKY gene was identified on  
16 Chr7, Chr9 and Chr12. The distribution of each SsWRKY genes was extremely irregular,  
17 indicating the reduction of the tandem duplication events in willow WRKY genes.

18 Gene clusters, defined as a single chromosome containing two or more genes [58], are very  
19 important for predicting co-expression genes or potential function of clustered genes in  
20 angiosperms [59]. According to this description, a total of 23 SsWRKY genes were clustered  
21 into 11 clusters in willow (Fig. 2). The chromosomal distribution of gene cluster was irregular,  
22 and only seven chromosomes were identified to have gene clusters. Three clusters, including  
23 seven SsWRKY genes, were found on Chr2, and two clusters were found on both Chr6 and  
24 Chr14. Only one cluster was distributed on each of Chr3, Chr8, Chr10 and Chr18, whereas  
25 none was identified on other eleven chromosomes. Further analysis of SsWRKY  
26 chromosomal distribution showed that a high WRKY gene density region in only 2.23 Mb  
27 regions on Chr2, which had also been observed in rice and poplar [7, 13].

# 1 Phylogenetic analysis and classification of WRKY genes in willow

2 In order to get a better separation of different groups and subgroups in SsWRKY genes, a  
3 total of 185 WRKY domains, including 82 AtWRKY domains and 103 SsWRKY domains,  
4 were used to construct the NJ phylogenetic tree. On the basis of the phylogenetic tree and  
5 structural features of WRKY domains, all 85 SsWRKY genes were clustered into three main  
6 groups (Fig. 3). Nineteen members containing two WRKY domains and C<sub>2</sub>H<sub>2</sub>-type zinc finger  
7 motifs were categorized into group I, except SsWRKY78, which contains only one WRKY  
8 domain and two zinc finger motifs. Domain acquisition and loss events appear to have shaped  
9 the WRKY family [60, 61]. Thus, SsWRKY78 may have evolved from a two-domain WRKY  
10 gene but lost one WRKY domain during evolution. Additionally, as shown in Fig. 3,  
11 SsWRKY78 shows high similarities to SsWRKY40N, implying a common origin of their  
12 domains. The similar phenomenon was also found in PtWRKY90 of poplar [7].

13 The largest number of SsWRKY genes, comprising a single WRKY domain and C<sub>2</sub>H<sub>2</sub> zinc  
14 finger motif, were categorized into group II. SsWRKY genes of group II could be further  
15 divided into five subgroups: IIa, IIb, IIc, IId and IId. As shown in Fig. 3, subgroup IIa (4  
16 members) and IIb (8 members) were clustered into one clade, as well as subgroup IId (13  
17 members) and IId (11 members). Strikingly, SsWRKY genes in subgroup IIc (21 members)  
18 and group IC are classified into one clade, suggesting that group II genes are not  
19 monophyletic and the group IIc WRKY genes may evolve from the group I genes by the loss  
20 of the WRKY domain in N-terminal. As shown in Fig. 3 and Fig. 4, SsWRKY23,  
21 SsWRKY34 and their orthologous genes, AtWRKY49, PtWRKY39, PtWRKY57,  
22 PtWRKY34 and PtWRKY32, seem to form a new subgroup, and shown to be closer to the  
23 group III according to the phylogenetic analysis. However, SsWRKY23 and SsWRKY34  
24 exhibit the zinc finger motif C-X<sub>4</sub>-C-X<sub>21</sub>-H-X-H and C-X<sub>4</sub>-C-X<sub>23</sub>-H-X-H as observed in the  
25 subgroup IIc and group IC. Thereby, they were classified into subgroup IIc in this study.

26 Different from the C<sub>2</sub>H<sub>2</sub> zinc finger pattern in group I and II, group III WRKY genes (7  
27 members), broadly considered as playing vital roles in plant evolution process and  
28 adaptability, contained one WRKY domain and a C-X<sub>7</sub>-C-X<sub>23</sub>-H-X-C zinc finger motif.

1 Intriguingly, a subgroup IIIb containing a  $CX_7CX_nHX_1C$  ( $n \geq 24$ ) zinc finger motif was  
2 identified in rice and barley [11, 13]. However, this  $C-X_7-C-X_n-H-X-C$  ( $n \geq 24$ ) zinc finger  
3 motif was never found in poplar, grape, *Arabidopsis* and willow, suggesting that this feature  
4 perhaps only belong to monocotyledonous species.

5 In order to obtain a better study in woody plant species, a phylogenetic tree based on the  
6 WRKY domains between willow and poplar was constructed (Fig. 4). The tree showed that  
7 most of the WRKY domains from willow and poplar were clustered into sister pairs,  
8 suggesting that gene duplication events played prominent roles in the evolution and expansion  
9 of WRKY gene family. Furthermore, a total of twenty SsWRKY domains show extremely the  
10 same domains (similarity: 100%) to poplar, i.e., SsWRKY39 and PtWRKY9, SsWRKY39  
11 and PtWRKY9, SsWRKY39 and PtWRKY9, SsWRKY39 and PtWRKY9, and so on. Further  
12 functional analyses of these genes in willow or poplar will provide a useful reference for  
13 another one.

## 14 **The ortholog of SsWRKY genes in *Arabidopsis* and poplar**

15 The clustering of orthologous genes emphasizes the conservation and divergence of gene  
16 families, and they may contain the same functions [9]. In this study, a phylogeny-based  
17 method was used to identify the putative orthologous SsWRKY genes in *Arabidopsis* and  
18 poplar (Fig. 3 and Fig. 4), and BLAST-based method (Bi-direction best hit) was used to  
19 confirm the true orthologs. The WRKY genes of group I contained two WRKY domains, and  
20 both of them were used to construct the phylogenetic trees. To avoid the mistakes of  
21 orthologous genes in group I, the members of group I WRKY genes were considered as  
22 orthologous genes unless the same phylogenetic relationship can be detected between  
23 N-termini and C-termini in the phylogenetic tree. For example, SsWRKY37 and AtWRKY44  
24 were considered as an orthologous gene pair because they clustered into a clade of their  
25 N-termini and C-termini (Fig. 3), while SsWRKY80 and PtWRKY30 were excluded from  
26 orthologous gene pairs due to their different clusters of N-termini and C-termini (Fig. 4).  
27 Totally, 75 orthologous gene pairs were found between willow and *Arabidopsis*, less than 82

1 orthologous gene pairs between willow and poplar (Table 1), which was congruent with the  
2 evolutionary relationship among the three plant species.

### 3 **Evolutionary analysis of WRKY III genes in willow**

4 The WRKY III genes were considered as the evolutionary youngest groups, and played  
5 crucial roles in the process of plant growth and resistance. In order to further probe the  
6 duplication and diversification of WRKY III genes after the divergence of the monocots and  
7 dicots, a phylogenetic tree was constructed using 65 WRKY III genes from *Arabidopsis* (13),  
8 rice (29), poplar (10), willow (7) and grape (6). As shown in Fig. 5, willow SsWRKY III  
9 genes were closer to the eurosids I group (poplar and grape) than eurosids II group  
10 (*Arabidopsis*) and monocots (rice). Meanwhile, most *Arabidopsis* and rice WRKY III genes  
11 formed the relatively independent clades, suggesting that two gene duplication events,  
12 including tandem and segmental duplication, perhaps were the main factors in the expansion  
13 of WRKY III genes in *Arabidopsis* and rice. What's more, the results also indicated that  
14 WRKY III might arise after the divergence of the *Arabidopsis* (eurosids I) and eurosids II  
15 (poplar, willow and grape). The study by Ling et al. in cucumber [9] showed the similar  
16 results and hence proved the validity. Interestingly, seven rice WRKY III genes (OsWRKY55,  
17 84, 18, 52, 46, 114 and 97) contained the variant domain WRKYGEK, but the variant was not  
18 found in other four dicots, implying that this may be a feature of WRKY III genes in  
19 monocots and these OsWRKY genes may respond to different environmental signals.

20 According to the comparison of the number of WRKY III genes in the five observed plants,  
21 the number is smaller in eurosids I (poplar, grape and willow) than *Arabidopsis* (eurosids II)  
22 and rice (monocots), which may be caused by different patterns of duplication events. Genes  
23 generated by duplication events are not stable, and can be retained or lost due to different  
24 selection pressure and evolution [62]. In order to determine which selection pressure played  
25 prominent roles in the expansion of willow WRKY III genes, we estimated the Ka/Ks ratios  
26 for all pairs (21 pairs) of willow WRKY III genes. As shown in Fig. 6, all the Ka/Ks ratios



1 were less than 0.5, suggesting willow WRKY III genes had mainly been subjected to strong  
2 purifying selection and they were slowing evolving at the protein level.

### 3 **Exon–intron structures of SsWRKY genes**

4 The exon-intron structures of multiple gene families play crucial roles during plant evolution.  
5 As shown in Fig. 7, the SsWRKY gene phylogenetic tree and the corresponding exon-intron  
6 structures are shown in A and B, respectively. Exon-intron structures of each group were  
7 shown in Fig. 7B, a large number of WRKY genes had two to five introns (94%, 80 of 85),  
8 including 8 WRKY genes contained one intron; 39 contained two introns; 13 contained three  
9 introns; 15 contained four introns and 5 contained five introns. The number of exons in  
10 remaining WRKY genes was quite different: SsWRKY49, SsWRKY76 and SsWRKY78 had  
11 six, eleven and ten introns, respectively; SsWRKY17 had the largest number of introns  
12 (seventeen introns), while no intron was found in SsWRKY12. The intron acquisition or loss  
13 occurred during the evolution of WRKY gene family, while WRKY genes in the same group  
14 shared the similar number of introns [6]. In our study, most of WRKY genes in group I had  
15 three to six introns, except SsWRKY76 and SsWRKY78, which might acquire some introns  
16 during evolution. The number of introns of WRKY genes in group II was extremely different,  
17 ranging from one to five introns, except SsWRKY17 with 17 introns and SsWRKY12 with  
18 zero intron might obtain or loss some introns during evolution. Strikingly, WRKY genes in  
19 group III had the most stable number of introns with all of seven WRKY III genes had two  
20 introns, suggesting that WRKY III genes may be the most stable genes in the environmental  
21 stress. The stable number of introns in SsWRKY III genes was consistent with the results of  
22 Ka/Ks analysis, which reflected that purifying selection pressure played vital roles in willow  
23 WRKY III genes.

24 A great deal of research in WRKY genes proved that nearly all of the WRKY genes  
25 contained an intron in their WRKY core domains [2, 6-9, 30]. According to the further  
26 analysis of SsWRKY genes, two major types of splicing introns, R-type and V-type, introns  
27 were observed in numerous SsWRKY domains. The R-type intron was spliced exactly at the



1 R residue, about five amino acids before the first Cys residue in the C<sub>2</sub>H<sub>2</sub> zinc finger motif.  
 2 The V-type intron was localized before the V residue, six amino acids after the second Cys  
 3 residue in the C<sub>2</sub>H<sub>2</sub> zinc finger motif. As shown in Fig. 7B, the R-type introns could be  
 4 observed in more groups, including group IC, subgroup IIc, IId, IIe and group III, while  
 5 V-type introns were only observed in subgroup IIa and IIb. Moreover, there was no intron  
 6 found in group IN. The similar results were also observed in *Arabidopsis*, poplar and rice,  
 7 suggesting that the special distribution of introns in WRKY domains was a feature of WRKY  
 8 family.

## 9 Identification of gene duplication events and conserved motifs in 10 willow

11 Gene duplication events were always considered as the vital sources of biological evolution  
 12 [63, 64]. Two or more adjacent homologous genes located on a single chromosome were  
 13 considered as tandem duplication events (TDs), while homologous gene pairs between  
 14 different chromosomes were defined as segmental duplication events (SDs) [10]. In our study,  
 15 a total of 33 homologous gene pairs, including 66 SsWRKY genes, were identified to  
 16 participate in gene duplication events. The composition of gene duplication events in each  
 17 group in ascending order was group I: 73.7% (14 of 19), group II: 78% (46 of 59) and group  
 18 III: 85.7% (6 of 7). Among the 33 homologous gene pairs, none of them appeared to have  
 19 undergone TDs, on the contrary, all of the 66 genes (77.6% of all SsWRKY genes)  
 20 participated in SDs, implying that segmental duplication events played major roles in the  
 21 expansion of willow WRKY genes.

22 WRKY genes shared more functional and homologies in their conserved WRKY core  
 23 domains (about 60 residues), while the rest sequences of WRKY genes shared a little [2]. In  
 24 order to get a more comprehensive understanding of the structural feature in WRKY domains,  
 25 the conserved motifs of SsWRKY genes were predicted using the online program MEME  
 26 (Fig. 8 and Table 2). Among the 20 putative motifs, motifs 1, 2, 3 and 5, broadly distributed  
 27 across SsWRKY genes, were characterized as the WRKY conserved domains. The motif 6

1 was characterized as nuclear localization signals (NLS), which mainly distributed in subgroup  
2 II d and IIe and group III. Some other motifs with poorly defined recently were also predicted  
3 by MEME: the motif 4 was only found in group IC and subgroup IIc; motifs 7 and 9 were  
4 limited to subgroup IIa and IIb; the motif 8 was found in group I and a few genes of subgroup  
5 IIc; motifs 10, 13, 15 and 17 were unique in subgroup IId; the motif 12 was only observed in  
6 subgroup IIb; the motif 16 was mainly found in group II; the motif 18 was found in subgroup  
7 IIc; motifs 19 and 20 were only observed in subgroup I. The distinct conserved motifs of  
8 different groups could be an important foundation for future structural and functional study in  
9 WRKY gene family.

10 Some other important motifs, including Leu zipper motif, HARF, LXXLL and LXLXLX,  
11 could be also identified in WRKY genes. Using the online program 2ZIP, the conserved Leu  
12 zipper motif, described as a common hypothetical structure to DNA binding proteins [65],  
13 was identified in only two SsWRKY genes (SsWRKY61 and SsWRKY39). With manual  
14 inspection, the conserved HARF (RTGHARFRR[A/G]P) motifs, whose putative functions  
15 were not distinguished clearly, were only observed in seven WRKY genes of subgroup IId,  
16 including SsWRKY82, 33, 45, 81, 9, 30 and 56. In the meantime, the conserved LXXLL and  
17 LXLXLX (L: Leucine; X: any amino acid) motifs, which respectively defined as the  
18 co-activator and active repressor motifs, were also found in SsWRKY genes. A total of seven  
19 SsWRKY genes (SsWRKY19, 45, 72, 61, 76, 30 and 59) contained the helical motif LXXLL,  
20 whereas eight genes (SsWRKY66, 26, 35, 81, 83, 75, 73 and 3) shared the LXLXLX motif.  
21 The plenty of conserved motifs in WRKY genes with different lengths and variant functions,  
22 suggesting that the WRKY genes might play more vital roles in gene regulatory network.

## 23 **Distinct expression profiles of SsWRKY genes in various tissues**

24 In order to gain more information about the roles of WRKY genes in willow, RNA-seq data  
25 from the sequenced genotype were used to quantify the expression level of WRKY genes in  
26 five tissues of *Salix suchowensis*. As illustrated in Fig. 9, the expression of all 85 SsWRKY  
27 genes were detected in at least one of the five examined tissues, such as 84 genes in roots, 80

in stems, 84 in barks, all in buds and 73 in leaves. Meanwhile, the cluster analysis of the expression pattern in five tissues showed that SsWRKY genes shared more similarities between stem and leaf, as well as bark and bud, and root was more similar to the clade formed by bark and bud. The results detected here were consistent with their biological characteristics. SsWRKY38, not detected in roots and leaves, was also lowly expressed in other tissues. Similarly, SsWRKY74, not detected in stems, barks and leaves, was only expressed in roots and buds with extremely low levels. Among the five genes not expressed in stems, SsWRKY66, 74 and 79 were also not detected in leaves. The largest number of expressed or unexpressed SsWRKY genes (12 genes) was found in buds or leaves, respectively, suggesting that WRKY genes might play more roles in buds than leaves.

According to the expression annotation of 85 SsWRKY genes by RPKM method in Fig. 9 and Table S1, the total transcript abundance of SsWRKY genes in tender root (RPKM = 1181.21), bark (RPKM = 1363.01) and vegetative bud (RPKM = 928.58) was relatively larger than that in other two tissues, including non-lignified stem (RPKM = 537.88) and young leaf (RPKM = 349.84). As shown in Table S1, SsWRKY81 (RPKM = 97.75), the most expressed SsWRKY genes in roots, was also expressed in other four tissues, though the expression levels were relatively low; SsWRKY56 (RPKM = 32.54), the most expressed SsWRKY genes in stem, was also highly expressed in other examined tissues. Similarly, SsWRKY67, the most expressed SsWRKY genes in barks (RPKM = 188.16), was also detected in vegetative buds (RPKM = 82.07) and young leaves (RPKM = 26.11) with high expression levels. Similarly, SsWRKY6 (RPKM = 26.31), the most expressed genes in leaves, was also highly detected in other tissues. A few genes, i.e., SsWRKY52, SsWRKY2 and SsWRKY35, were expressed highly in barks, but lowly in other four tissues. The results mentioned above may be an important foundation for the specific expression analysis of each WRKY gene in willow.

# 1 Discussion

2 WRKY genes are the induced plant TFs, which can specifically interact with the W-box to  
3 regulate the expressions of downstream target genes. They also play prominent roles in  
4 diverse physiological and growing processes, especially in various abiotic and biotic stress  
5 responses in plants. Previous studies about the features and functions of WRKY family have  
6 been conducted in many model plants, including *Arabidopsis* for annual herbaceous dicots [2],  
7 grape for perennial dicots [6], poplar for woody plants and rice for monocots [7, 13]. Here,  
8 the comprehensive analyses of WRKY family in willow (*Salix suchowensis*) would not only  
9 provide valuable information for future functional analysis of WRKY genes in woody plants,  
10 but also provide an important reference to investigate the complex structures, evolution and  
11 gene expansion in this gene superfamily. In this study, a total of 85 SsWRKY genes were  
12 identified from willow, accompanying with analyses of their complex structures,  
13 classification, gene expansion patterns, conserved motifs and distinct expression profiles.

14 Comparing the two phylogenetic trees based on the SsWRKY domains (Fig. 3) and  
15 proteins (Fig. 7 A), we obtained the nearly same classification of all SsWRKY genes,  
16 suggesting that the conserved WRKY domain is an indispensable unit in WRKY genes. The  
17 variation of the WRKYGQK heptapeptide may influence the proper DNA-binding ability of  
18 WRKY genes [17, 18]. A recent binding study by Brand et al. disclosed that a reciprocal Q/K  
19 change of the WRKYGQK heptapeptide might result in different DNA-binding specificities  
20 of the respective WRKY genes [56]. For instance, the soybean WRKY genes, GmWRKY6  
21 and GmWRKY21, which contains the WRKYGKK variant, can't bind normally to the W-box  
22 ([C/T]TGAC[T/C]) [66]. Another NtWRKY12 gene in tobacco with the WRKYGKK variant  
23 recognizes another binding sequence 'TTTTCCAC' instead of normal W-box [67]. Strikingly,  
24 many WRKY genes with WRKYGKK variant recognize a much more degenerate consensus  
25 with only a central GAC-core motif, i.e., AtWRKY50 in *Arabidopsis* [56]. Therefore, further  
26 investigation of the functions and binding specificities of the variants of WRKYGQK  
27 heptapeptide in plants would be very interesting. In our study, four WRKY genes

(SsWRKY14, SsWRKY23, SsWRKY38 and SsWRKY78) had a single mismatched amino acid in their conserved WRKYGQK heptapeptide (Fig. 1), including WRKYGKK, WKKYGQK and WRKYGRK. The variants detected in willow were extremely congruent with that in another salicaceous plant, poplar, which also contains the three variants in seven PtWRKY genes [7]. Additionally, two variants, WRKYGKK and WRRKGQK, were found in grape and tomato [6, 8]; WRKYGKK, the most common variant in plants, was the only one found in castor bean and cucumber [9, 68]. The variants may be different between dicots and monocots. Four variants, including WQKYGQK, WRKYGKK, WSKYGQM and WRKYGEK, were found in barley [11]. Meanwhile, the largest number of variants was found in rice [13], including WQKYGQK, WRKYGEK, WIKYGQK, WRKYSEK, WKKYGQK, WKRYGQK, WSKYEQK and WRKYGKK, perhaps due to its various habitats. Strikingly, WRKYGEK, a prevalent variant in plants, was only found in WRKY III genes of rice and barley among the above plants examined, implying that this variant may be a feature of WRKY III genes in monocots and they may respond to different environmental signals. Moreover, many previous studies have disclosed that the binding specificities of variable WRKYGQK heptapeptide vary tremendously [56], however, few studies were shown about the effect of variable zinc finger motif. In this study, four WRKY domains (SsWRKY76C, SsWRKY64, SsWRKY12 and SsWRKY28) without complete zinc finger motif may lack the ability of interacting with W-box, as well as PtWRKY83, 40, 95 and 10 in poplar [7]. It is still indispensable to further investigate the function or the expression patterns of the regulated gene targets in the variant sequences of the WRKY conserved domains.

Different classification methods may lead to different numbers of WRKY genes in each group. The classification method in our study was categorized as described in *Arabidopsis*, grape, cucumber, castor bean and many other plant species [2, 6, 9, 68]. According to this method, the WRKY genes were classified into three main groups (I, II and III), with five subgroups in group II (IIa, IIb, IIc, IId and IIe) based on the number of WRKY domains and the features of diverse zinc finger motifs. However, the strategy described in rice and poplar was a little different [7, 13], and they classified the subgroup IIc categorized above into a new

subgroup Ib based on the fact that the C-termini of group I and the domains of the above subgroup IIc shared more similar consensus structures. At the meantime, subgroup IId and IIe categorized above were reclassified into subgroup IIc and IId, respectively. Thus, the number of WRKY genes in poplar and rice was different from other plant species (Table 3). With the same classification method as described in *Arabidopsis* and many other plants, the number of different groups in poplar was as follows: group I: 23, subgroup IIa: 5, IIb: 9, IIc: 31, IId: 13, IIe: 13 and group III: 10, and the number of OsWRKY genes in rice: group I: 14, subgroup IIa: 4, IIb: 8, IIc: 20, IId: 7, IIe: 11 and group III: 36. WRKY genes of subgroup IIa, the smallest number of members, appear to play crucial roles in regulating stress responses (both biotic and abiotic) [3]. As illustrated in Table 3, the WRKY genes of subgroup IIa and IIb in willow are extremely similar to that of other plant species, suggesting that all SsWRKY genes of these subgroups have been identified. Subgroup IIa genes, the smallest number of members, appear to play many important roles in regulating biotic and abiotic stress responses [3]. Nevertheless, the number of WRKY III in eurosids I group, such as cucumber (6), poplar (10), grape (6) and willow (7) is less than that of eurosids II (*Arabidopsis*: 14) and monocots (rice: 36), suggesting that different duplication events or selection pressures occurred in WRKY III genes after the divergence of eurosids I and eurosids II group. Interestingly, the previous study in *Arabidopsis* showed that nearly all WRKY III members respond to diverse biotic stresses, suggesting that this group probably evolved with the increasing biological requirements and the larger number of WRKY III genes in *Arabidopsis* and rice probably due to their various biotic stresses during evolution.

WRKY transcription factors play important roles in the regulation of developmental processes and response to biotic and abiotic stress [56]. The evolutionary relationship of WRKY gene family promises to obtain significant insights into how biotic and abiotic stress responses from single cellular aquatic algae to multicellular flowering plants [57]. The first work by Eulgem et al. defined the seven major groups of WRKY genes observed in flowering plants, which has proven over time to be an accurate representation of groups of WRKY genes [2, 3]. Previous studies hypothesized that group I WRKY genes were generated by

domain duplication of a proto-WRKY gene with a single WRKY domain, group II WRKY genes evolved through the subsequent loss of N-terminal WRKY domain, and group III genes evolved from the replacement of conserved His residue with a Cys residue in zinc motif [13]. However, recent study proposed two alternative hypotheses of WRKY gene evolution [57]: the "Group I Hypothesis" suggests that all WRKY genes in higher plants evolved from group I genes, while the "IIa + b Separate Hypothesis" considers that subgroup IIa and IIb with their hallmark V-type intron are evolved from a single domain of ancestral algal WRKY gene instead of evolving from group I genes. Additionally, another recent study by Brand et al. concluded that subgroup IIc WRKY genes evolved directly from IIc-like ancestral WRKY domains, and group I genes evolved independently due to a duplication of the IIc-like ancestral WRKY domains [56]. In his study, subgroup IIa genes evolved from group I genes through loss of their N-terminal domains; subgroup IIb genes were descendants from IIa genes, because IIb representatives can only be found in monocots and dicots; subgroup IId genes evolved most probably from IIa, and IIe are most likely the descendants from IId WRKY genes; and group III WRKY genes are considered as the evolutionary youngest genes. Phylogenetic analysis in our study shows that subgroup IIc and group IC are evolutionarily close, as well as subgroups IIa and IIb, subgroups IId and IIe, and this result is consistent with the conclusion drew by Brand et al [56]. Additionally, the V-type introns of SsWRKY genes are only found in subgroup IIa and IIb, while R-type introns are found in other groups except group IN. The results are congruent with the "IIa + b Separate Hypothesis". Therefore, further information is still required to determine the accurate evolutionary relationship of WRKY gene family.

Gene duplication events played prominent roles in a succession of genomic rearrangements and expansions, and it is also the main motivation of plants evolution [69]. The gene family expansion occurs via three mechanisms: tandem duplication events (TDs), segmental duplication events (SDs) and transposition events [70], and we only focused on the tandem and segmental duplication events in this study. In willow, a total of 66 SsWRKY genes were identified to participate in gene duplication events, and all of these genes appeared to have



undergone SDs. In poplar, only one homologous gene pair participated in TDs, while 29 of 42 (69%) homologous gene pairs were determined to participate in SDs. The WRKY gene expansion patterns in willow and poplar perhaps showed that SDs were the main factors in the expansion of WRKY genes in woody plants. However, in cucumber, no gene duplication events have occurred in CsWRKY gene evolution, probably because there were no recent whole-genome duplication and tandem duplication in cucumber genome [71]. In rice and *Arabidopsis*, many WRKY genes were generated by TDs, which was incongruent with the duplication events in willow, poplar and cucumber. The different WRKY gene expansion patterns of the above plant species could be due to their different life habits and selection pressures in a large scale.

The WRKY gene family plays crucial roles in response to biotic and abiotic stresses, as well as diverse physiological and developmental processes in plant species. Because of the lack of researches on the function of willow WRKY genes, our study provided putative functions of SsWRKY genes by comparing the orthologous genes between willow and *Arabidopsis*. The details of the functions or regulations of AtWRKY genes can be obtained from TAIR (<http://www.arabidopsis.org/>). For example, AtWRKY2, the ortholog to SsWRKY6, which highly expressed in the five examined tissues, plays important roles in seed germination and post germination growth [72]. AtWRKY33, the ortholog to SsWRKY1, 35, 55 and 84, influences the tolerance to NaCl, inc sensitivity to oxidative stress and abscisic acid [25]. A large number of AtWRKY genes, i.e. AtWRKY3, 4, 18, 53, 41, work in the resistance to *Pseudomonas syringae* [73-76], so do their orthologs in willow (SsWRKY42, 47, 39, 79, 20 and 70). Based on the comparison of willow WRKY genes with their *Arabidopsis* orthologs, we could speculate that the functional divergence of SsWRKY genes has played prominent roles in the responses to various stresses.



# Conflict of Interests

The authors declare that there is no conflict of interests regarding the publication of this paper.

# Acknowledgment

We thank the Fundamental Research Funds for the Central Non-profit Research Institution of CAF (CAFYBB2014QB015), National Basic Research Program of China (973 Program) (2012CB114505) and the National Natural Science Foundation of China (31570662, 31500533 and 61401214). We also acknowledge supports from Key Projects in the National Science & Technology Pillar Program during the Twelfth Five-year Plan Period (NO.2012BAD01B07), and Natural Science Foundation of the Jiangsu Higher Education Institutions (14KJB520018). This work is also enabled by the Innovative Research Team Program of the Educational Department of China, the Innovative Research Team Program in Universities of Jiangsu Province, Scientific Research Foundation for Advanced Talents and Returned Overseas Scholars of Nanjing Forestry University and the PAPD (Priority Academic Program Development) program at Nanjing Forestry University.

# References

- Jang, J.-Y., C.-H. Choi, and D.-J. Hwang, The WRKY Superfamily of Rice Transcription Factors. *The Plant Pathology Journal*, 2010. **26**(2): p. 110-114.
- Eulgem, T., The WRKY superfamily of plant transcription factors. *Trends in Plant Science*, 2000. **5**(5): p. 199-206.
- Rushton, P.J., et al., WRKY transcription factors. *Trends Plant Sci*, 2010. **15**(5): p. 247-58.
- Ishiguro, S. and K. Nakamura, Characterization of a cDNA encoding a novel DNA-binding protein, SPF1, that recognizes SP8 sequences in the 5' upstream regions of genes coding for sporamin and  $\beta$ -amylase from sweet potato. *MGG Molecular & General Genetics*, 1994. **244**(6).
- Giacomelli, J.I., et al., Expression analyses indicate the involvement of sunflower WRKY transcription factors in stress responses, and phylogenetic reconstructions

- 1 reveal the existence of a novel clade in the Asteraceae. Plant Science, 2010. **178**(4): p.
- 2 398-410.
- 3 6. Guo, C., et al., Evolution and expression analysis of the grape (*Vitis vinifera* L.)
- 4 WRKY gene family. J Exp Bot, 2014. **65**(6): p. 1513-28.
- 5 7. He, H., et al., Genome-wide survey and characterization of the WRKY gene family in
- 6 *Populus trichocarpa*. Plant Cell Rep, 2012. **31**(7): p. 1199-217.
- 7 8. Huang, S., et al., Genome-wide analysis of WRKY transcription factors in *Solanum*
- 8 *lycopersicum*. Mol Genet Genomics, 2012. **287**(6): p. 495-513.
- 9 9. Ling, J., et al., Genome-wide analysis of WRKY gene family in *Cucumis sativus*.
- 10 BMC Genomics, 2011. **12**: p. 471.
- 11 10. Liu, J.J. and A.K. Ekramoddoullah, Identification and characterization of the WRKY
- 12 transcription factor family in *Pinus monticola*. Genome, 2009. **52**(1): p. 77-88.
- 13 11. Mangelsen, E., et al., Phylogenetic and comparative gene expression analysis of
- 14 barley (*Hordeum vulgare*) WRKY transcription factor family reveals putatively
- 15 retained functions between monocots and dicots. BMC Genomics, 2008. **9**: p. 194.
- 16 12. Pnueli, L., et al., Molecular and biochemical mechanisms associated with dormancy
- 17 and drought tolerance in the desert legume *Retama raetam*. The Plant Journal, 2002.
- 18 **31**(3): p. 319-330.
- 19 13. Wu, K.L., The WRKY Family of Transcription Factors in Rice and *Arabidopsis* and
- 20 Their Origins. DNA Research, 2005. **12**(1): p. 9-26.
- 21 14. Xu, X., et al., Physical and functional interactions between pathogen-induced
- 22 *Arabidopsis* WRKY18, WRKY40, and WRKY60 transcription factors. Plant Cell,
- 23 2006. **18**(5): p. 1310-26.
- 24 15. Ciolkowski, I., et al., Studies on DNA-binding selectivity of WRKY transcription
- 25 factors lend structural clues into WRKY-domain function. Plant Mol Biol, 2008.
- 26 **68**(1-2): p. 81-92.
- 27 16. Sun, C., A Novel WRKY Transcription Factor, SUSIBA2, Participates in Sugar
- 28 Signaling in Barley by Binding to the Sugar-Responsive Elements of the iso1
- 29 Promoter. The Plant Cell Online, 2003. **15**(9): p. 2076-2092.
- 30 17. Duan, M.R., et al., DNA binding mechanism revealed by high resolution crystal
- 31 structure of *Arabidopsis thaliana* WRKY1 protein. Nucleic Acids Res, 2007. **35**(4):
- 32 p. 1145-54.
- 33 18. Maeo, K., et al., Role of conserved residues of the WRKY domain in the
- 34 DNA-binding of tobacco WRKY family proteins. Biosci Biotechnol Biochem, 2001.
- 35 **65**(11): p. 2428-36.
- 36 19. Yamasaki, K., et al., Solution structure of an Arabidopsis WRKY DNA binding
- 37 domain. Plant Cell, 2005. **17**(3): p. 944-56.
- 38 20. Cormack, R.S., et al., Leucine zipper-containing WRKY proteins widen the spectrum
- 39 of immediate early elicitor-induced WRKY transcription factors in parsley.
- 40 Biochimica et Biophysica Acta (BBA) - Gene Structure and Expression, 2002.
- 41 **1576**(1-2): p. 92-100.

- 1 21. Ulker, B. and I.E. Somssich, WRKY transcription factors: from DNA binding  
2 towards biological function. *Curr Opin Plant Biol*, 2004. **7**(5): p. 491-8.
- 3 22. Ramamoorthy, R., et al., A comprehensive transcriptional profiling of the WRKY  
4 gene family in rice under various abiotic and phytohormone treatments. *Plant Cell*  
5 *Physiol*, 2008. **49**(6): p. 865-79.
- 6 23. Dong, J., C. Chen, and Z. Chen, Expression profiles of the *Arabidopsis* WRKY gene  
7 superfamily during plant defense response. *Plant Molecular Biology*, 2003. **51**(1): p.  
8 21-37.
- 9 24. Li, J., et al., WRKY70 modulates the selection of signaling pathways in plant  
10 defense. *Plant J*, 2006. **46**(3): p. 477-91.
- 11 25. Jiang, Y. and M.K. Deyholos, Functional characterization of *Arabidopsis*  
12 NaCl-inducible WRKY25 and WRKY33 transcription factors in abiotic stresses.  
13 *Plant Mol Biol*, 2009. **69**(1-2): p. 91-105.
- 14 26. Lagace, M. and D.P. Matton, Characterization of a WRKY transcription factor  
15 expressed in late torpedo-stage embryos of *Solanum chacoense*. *Planta*, 2004. **219**(1):  
16 p. 185-9.
- 17 27. Robatzek, S. and I.E. Somssich, Targets of AtWRKY6 regulation during plant  
18 senescence and pathogen defense. *Genes Dev*, 2002. **16**(9): p. 1139-49.
- 19 28. Johnson, C.S., TRANSPARENT TESTA GLABRA2, a Trichome and Seed Coat  
20 Development Gene of *Arabidopsis*, Encodes a WRKY Transcription Factor. *The*  
21 *Plant Cell Online*, 2002. **14**(6): p. 1359-1375.
- 22 29. Zhang, Z.L., et al., A rice WRKY gene encodes a transcriptional repressor of the  
23 gibberellin signaling pathway in aleurone cells. *Plant Physiol*, 2004. **134**(4): p.  
24 1500-13.
- 25 30. Zou, X., et al., A WRKY gene from creosote bush encodes an activator of the abscisic  
26 acid signaling pathway. *J Biol Chem*, 2004. **279**(53): p. 55770-9.
- 27 31. Du, L. and Z. Chen, Identification of genes encoding receptor-like protein kinases as  
28 possible targets of pathogen- and salicylic acid-induced WRKY DNA-binding  
29 proteins in *Arabidopsis*. *The Plant Journal*, 2008. **24**(6): p. 837-847.
- 30 32. Ding, M., et al., Genome-wide investigation and transcriptome analysis of the WRKY  
31 gene family in *Gossypium*. *Mol Genet Genomics*, 2015. **290**(1): p. 151-71.
- 32 33. Muthamilarasan, M., et al., Global analysis of WRKY transcription factor  
33 superfamily in *Setaria* identifies potential candidates involved in abiotic stress  
34 signaling. *Front Plant Sci*, 2015. **6**: p. 910.
- 35 34. Xiong, W., et al., Genome-wide analysis of the WRKY gene family in physic nut  
36 (*Jatropha curcas* L.). *Gene*, 2013. **524**(2): p. 124-32.
- 37 35. Zhang, Y. and L. Wang, The WRKY transcription factor superfamily: its origin in  
38 eukaryotes and expansion in plants. *BMC Evol Biol*, 2005. **5**: p. 1.
- 39 36. Santos, C.S., et al., Searching for resistance genes to *Bursaphelenchus xylophilus*  
40 using high throughput screening. *BMC Genomics*, 2012. **13**: p. 599.
- 41 37. Qiu, Y., Cloning and analysis of expression profile of 13 WRKY genes in rice.  
42 *Chinese Science Bulletin*, 2004. **49**(20): p. 2159.

- 1 38. Xie, Z., et al., Annotations and functional analyses of the rice WRKY gene  
2 superfamily reveal positive and negative regulators of abscisic acid signaling in  
3 aleurone cells. *Plant Physiol*, 2005. **137**(1): p. 176-89.
- 4 39. Dai, X., et al., The willow genome and divergent evolution from poplar after the  
5 common genome duplication. *Cell Res*, 2014. **24**(10): p. 1274-7.
- 6 40. Punta, M., et al., The Pfam protein families database. *Nucleic Acids Res*, 2012.  
7 **40**(Database issue): p. D290-301.
- 8 41. Camacho, C., et al., BLAST+: architecture and applications. *BMC Bioinformatics*,  
9 2009. **10**: p. 421.
- 10 42. Eddy, S.R., Profile hidden Markov models. *Bioinformatics*, 1998. **14**(9): p. 755-763.
- 11 43. Letunic, I., T. Doerks, and P. Bork, SMART: recent updates, new developments and  
12 status in 2015. *Nucleic Acids Res*, 2015. **43**(Database issue): p. D257-60.
- 13 44. Larkin, M.A., et al., Clustal W and Clustal X version 2.0. *Bioinformatics*, 2007.  
14 **23**(21): p. 2947-8.
- 15 45. Tamura, K., et al., MEGA6: Molecular Evolutionary Genetics Analysis version 6.0.  
16 *Mol Biol Evol*, 2013. **30**(12): p. 2725-9.
- 17 46. Chen, F., et al., Assessing performance of orthology detection strategies applied to  
18 eukaryotic genomes. *PLoS One*, 2007. **2**(4): p. e383.
- 19 47. Suyama, M., D. Torrents, and P. Bork, PAL2NAL: robust conversion of protein  
20 sequence alignments into the corresponding codon alignments. *Nucleic Acids Res*,  
21 2006. **34**(Web Server issue): p. W609-12.
- 22 48. Yang, Z., PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol*,  
23 2007. **24**(8): p. 1586-91.
- 24 49. Hu, B., et al., GSDS 2.0: an upgraded gene feature visualization server.  
25 *Bioinformatics*, 2015. **31**(8): p. 1296-7.
- 26 50. Gu, Z.L., et al., Extent of gene duplication in the genomes of Drosophila, nematode,  
27 and yeast. *Molecular Biology and Evolution*, 2002. **19**(3): p. 256-262.
- 28 51. Bailey, T.L., et al., MEME: discovering and analyzing DNA and protein sequence  
29 motifs. *Nucleic Acids Res*, 2006. **34**(Web Server issue): p. W369-73.
- 30 52. Bornberg-Bauer, E., Computational approaches to identify leucine zippers. *Nucleic*  
31 *Acids Research*, 1998. **26**(11): p. 2740-2746.
- 32 53. Li, H. and R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler  
33 transform. *Bioinformatics*, 2009. **25**(14): p. 1754-60.
- 34 54. Wagner, G.P., K. Kin, and V.J. Lynch, Measurement of mRNA abundance using  
35 RNA-seq data: RPKM measure is inconsistent among samples. *Theory Biosci*, 2012.  
36 **131**(4): p. 281-5.
- 37 55. Gentleman, R.C., et al., Bioconductor: open software development for computational  
38 biology and bioinformatics. *Genome Biol*, 2004. **5**(10): p. R80.
- 39 56. Brand, L.H., et al., Elucidating the evolutionary conserved DNA-binding specificities  
40 of WRKY transcription factors by molecular dynamics and in vitro binding assays.  
41 *Nucleic Acids Res*, 2013. **41**(21): p. 9764-78.

- 1 57. Rinerson, C.I., et al., The evolution of WRKY transcription factors. BMC Plant Biol,  
2 2015. **15**: p. 66.
- 3 58. Holub, E.B., The arms race is ancient history in *Arabidopsis*, the wildflower. Nat Rev  
4 Genet, 2001. **2**(7): p. 516-27.
- 5 59. Overbeek, R., et al., The use of gene clusters to infer functional coupling.  
6 Proceedings of the National Academy of Sciences, 1999. **96**(6): p. 2896-2901.
- 7 60. Ross, C.A., Y. Liu, and Q.J. Shen, The WRKY Gene Family in Rice (*Oryza sativa*).  
8 Journal of Integrative Plant Biology, 2007. **49**(6): p. 827-842.
- 9 61. Rossberg, M., et al., Comparative Sequence Analysis Reveals Extensive  
10 Microcolinearity in the Lateral Suppressor Regions of the Tomato, Arabidopsis, and  
11 Capsella Genomes. The Plant Cell, 2001. **13**(4): p. 979.
- 12 62. Zhang, J., Evolution by gene duplication: an update. Trends in Ecology & Evolution,  
13 2003. **18**(6): p. 292-298.
- 14 63. Chothia, C., et al., Evolution of the protein repertoire. Science, 2003. **300**(5626): p.  
15 1701-3.
- 16 64. Ohno, S., U. Wolf, and N.B. Atkin, Evolution from Fish to Mammals by Gene  
17 Duplication. Hereditas, 2009. **59**(1): p. 169-187.
- 18 65. McInerney, E.M., et al., Determinants of coactivator LXXLL motif specificity in  
19 nuclear receptor transcriptional activation. Genes & Development, 1998. **12**(21): p.  
20 3357-3368.
- 21 66. Zhou, Q.Y., et al., Soybean WRKY-type transcription factor genes, GmWRKY13,  
22 GmWRKY21, and GmWRKY54, confer differential tolerance to abiotic stresses in  
23 transgenic *Arabidopsis* plants. Plant Biotechnol J, 2008. **6**(5): p. 486-503.
- 24 67. van Verk, M.C., et al., A Novel WRKY transcription factor is required for induction  
25 of PR-1a gene expression by salicylic acid and bacterial elicitors. Plant Physiol, 2008.  
26 **146**(4): p. 1983-95.
- 27 68. Zou, Z., et al., Gene Structures, Evolution and Transcriptional Profiling of the WRKY  
28 Gene Family in Castor Bean (*Ricinus communis* L.). PLoS One, 2016. **11**(2): p.  
29 e0148243.
- 30 69. Vision, T.J., D.G. Brown, and S.D. Tanksley, The Origins of Genomic Duplications  
31 in *Arabidopsis*. Science, 2000. **290**(5499): p. 2114-2117.
- 32 70. Maher, C., L. Stein, and D. Ware, Evolution of *Arabidopsis* microRNA families  
33 through duplication events. Genome Res, 2006. **16**(4): p. 510-9.
- 34 71. Huang, S., et al., The genome of the cucumber, *Cucumis sativus* L. Nat Genet, 2009.  
35 **41**(12): p. 1275-81.
- 36 72. Jiang, W. and D. Yu, Arabidopsis WRKY2 transcription factor mediates seed  
37 germination and postgermination arrest of development by abscisic acid. BMC Plant  
38 Biol, 2009. **9**: p. 96.
- 39 73. Chen, C. and Z. Chen, Potentiation of developmentally regulated plant defense  
40 response by AtWRKY18, a pathogen-induced *Arabidopsis* transcription factor. Plant  
41 Physiol, 2002. **129**(2): p. 706-16.

- 1 74. Higashi, K., et al., Modulation of defense signal transduction by flagellin-induced  
2 WRKY41 transcription factor in *Arabidopsis thaliana*. Mol Genet Genomics, 2008.  
3 **279**(3): p. 303-12.
- 4 75. Lai, Z., et al., Roles of Arabidopsis WRKY3 and WRKY4 transcription factors in  
5 plant responses to pathogens. BMC Plant Biol, 2008. **8**: p. 68.
- 6 76. Murray, S.L., et al., Basal resistance against *Pseudomonas syringae* in *Arabidopsis*  
7 involves WRKY53 and a protein with homology to a nematode resistance protein.  
8 Mol Plant Microbe Interact, 2007. **20**(11): p. 1431-8.

9

## Figure 1(on next page)

### Comparison of the WRKY domain sequences from 85 SsWRKY genes.

The WRKY gene with the suffix -N and -C indicates the N-terminal and C-terminal WRKY domain of group I members, respectively. "-" has been inserted for the optimal alignment. Red indicates the highly conserved WRKYGQK heptapeptide, and the zinc finger motifs are highlighted in green. The position of a conserved intron is indicated by an arrowhead.



[illegible]

**Group I C**

SsWRKY55C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TFQG--	FWKKEVERASHD	LRAW	ITYEGCN	LDV
SsWRKY84C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TYQG--	FWKKEVERASHD	LRAW	ITYEGCN	LDV
SsWRKY1C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TFVG--	FWKKEVERASQD	LRAW	ITYEGCN	LDV
SsWRKY35C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSVG--	FWKKEVERASHD	LRAW	ITYEGCN	LDV
SsWRKY6C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	LKSV	ITYEGCN	LDV
SsWRKY51C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERAWD	LKSV	ITYEGCN	LDV
SsWRKY54C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	LKYV	ITYEGCN	LDV
SsWRKY4C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	LKPV	ITYEGCN	LDV
SsWRKY49C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	LKPV	ITYEGCN	LDV
SsWRKY42C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TTGP--	FWKKEVERAAAD	PRAW	ITYE	KNI
SsWRKY44C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TTAA--	FWKKEVERAAAD	PEAW	ITYEGCN	LDV
SsWRKY76C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERAPAD	PKAV	ITYEGCN	LDV
SsWRKY7C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERAPAD	PKAV	ITYEGCN	LDV
SsWRKY37C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERAPAD	PKAV	ITYEGCN	LDV
SsWRKY16C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERAYND	PKSV	ITYEGCN	LDV
SsWRKY65C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	SPRG--	FWKKEVERASHD	PKSV	ITYEGCN	LDV
SsWRKY40C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	PKSV	ITYEGCN	LDV
SsWRKY80C	LDDGVR	WRKYGGK	VGNGFN	PSRYYY	TSAG--	FWKKEVERASHD	PKSV	ITYEGCN	LDV

**Group II a**

SsWRKY22  
SsWRKY68  
SsWRKY39  
SsWRKY79

RDGVC **MRKYGK**VT--DNPCPRAYFK--FAP--S **SPVKKK**VSRIDQSVVATYEGG--NPH  
RDGVC **MRKYGK**VT--DNPCPRAYFK--FAP--S **SPVKKK**VSRIDQSVVATYEGG--NPH  
RDGVC **MRKYGK**VT--DNPCPRAYFK--FAP--S **SPVKKK**VSRIDQSVVATYEGG--NPH  
RDGVC **MRKYGK**VT--DNPCPRAYFK--SSP--S **SPVKKK**VSAENPTIVATYEGG--NAS

**Group II b**

SsWRKY17  
SsWRKY48  
SsWRKY24  
SsWRKY61  
SsWRKY66  
SsWRKY75  
SsWRKY73  
SsWRKY64

ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAA-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAV-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAV-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAG-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAA-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TVAP-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TVAP-G P--VRCABDRTITLTYYEGN PL  
ISDGGCC**WRKYGK**MAAGNFCPRARYR TMAA-G P--VRCABDRTITLTYYEGN PL

**Group II c**

```

LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY74LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY71LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY69LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY52LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY67LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY12LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY59LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY3LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY31LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY8LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY43LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY46LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY15LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY62LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY44LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY57LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY29LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY41LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY38LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
SsWRKY23LDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI
ADDDGKVRKRYKGGVNNKNSPFRSYYK YQG--LVKKRVLQLTKEGVVWITYEGNPI

```

[illegible][illegible]

**Group III**

↓

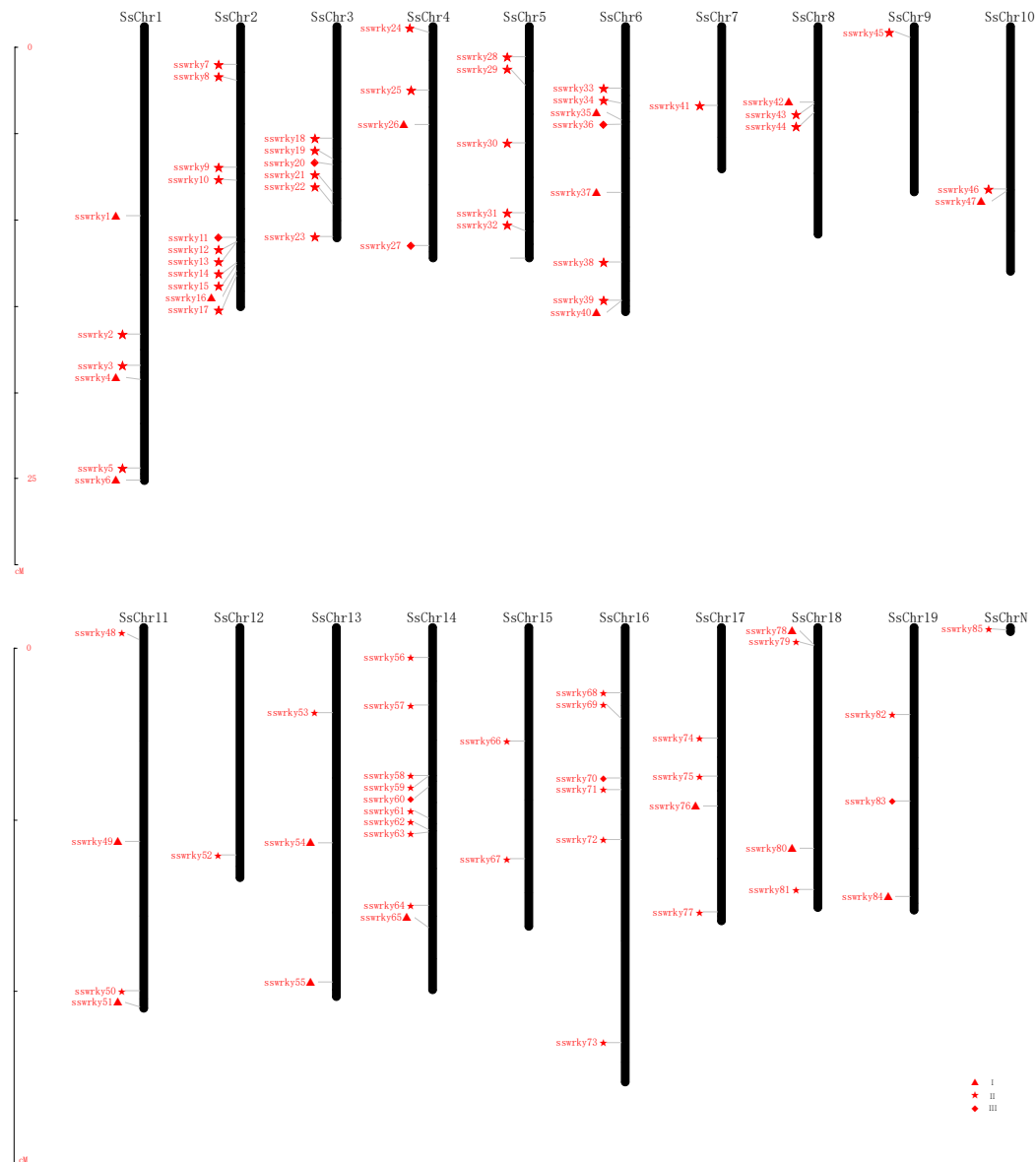
SsWRKY11 LDDG**Y** **MR**KY**G** **Q** LGAN **P**FR**Y** **Y**YR**Y** **Y**HR**S**Q**G** LAT **A** **W** **R** **D** **E** **N** **S** **H** **F** **E** **V** **N** **Y** **G** **E** **T** **C** **S** **Q**  
SsWRKY60 LDDG**Y** **MR**KY**G** **Q** LGAN **P**FR**Y** **Y**YR**Y** **Y**HR**S**Q**G** LAT **A** **W** **R** **D** **E** **N** **S** **H** **F** **E** **V** **N** **Y** **G** **E** **T** **C** **S** **Q**  
SsWRKY10 LDDG**Y** **MR**KY**G** **Q** LGAN **P**FR**Y** **Y**YR**Y** **Y**HR**S**Q**G** LAT **A** **W** **R** **D** **E** **N** **S** **H** **F** **E** **V** **N** **Y** **G** **E** **T** **C** **S** **Q**  
SsWRKY70 YDDG**Y** **MR**KY**G** **Q** LGT **P**FR**Y** **Y**YR**Y** **Y**HR**S**Q**G** LAT **A** **W** **R** **D** **E** **N** **S** **H** **F** **E** **V** **N** **Y** **G** **E** **T** **C** **S** **Q**  
SsWRKY93 P **D** **G** **T** **MR**KY**G** **Q** LG**S** **P**FR**Y** **Y**YR**Y** **Y**HR**S**Q**G** LAT **A** **W** **R** **D** **E** **N** **S** **H** **F** **E** **V** **N** **Y** **G** **E** **T** **C** **S** **Q**  
SsWRKY27 TDDG**Y** **MR**KY**G** **Q** **Y** **L** **N** **A** **K** **P** **R** **N** **R** **C** **T** **H** **K** **Y** **D** **D** **Q** **A** **T** **A** **T** **Q** **E** **P** **O** **T** **R** **Y** **V** **Y** **G** **H** **T** **C** **K** **N**  
SsWRKY36 TDDG**Y** **MR**KY**G** **Q** **Y** **L** **N** **A** **K** **P** **R** **N** **R** **C** **T** **H** **K** **Y** **D** **D** **Q** **A** **T** **A** **T** **Q** **E** **P** **O** **T** **R** **Y** **V** **Y** **G** **H** **T** **C** **K** **N**



## Figure 2(on next page)

### Chromosomal location of SsWRKY genes.

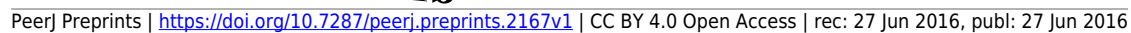
Red triangle indicates group I, red star indicates group II and red diamond indicates group III. The chromosome numbers are given at the top of each chromosome and the left side of each chromosome is related to the approximate physical location of each WRKY gene. Only one unmapped SsWRKY gene is shown on SsChrN.



# Figure 3(on next page)

## Phylogenetic tree of WRKY domains from willow and *Arabidopsis*.

The phylogenetic tree was constructed using the neighbor-joining method in MEGA 6.0. The WRKY genes with the suffix 'N' and 'C' indicate the N-terminal and the C-terminal WRKY domains of group I, respectively. The different colors indicate different groups (I, II and III) or subgroups (IIa, b, c, d and e) of WRKY domains. Circles indicate WRKY genes from willow, and diamonds represent genes from *Arabidopsis*.

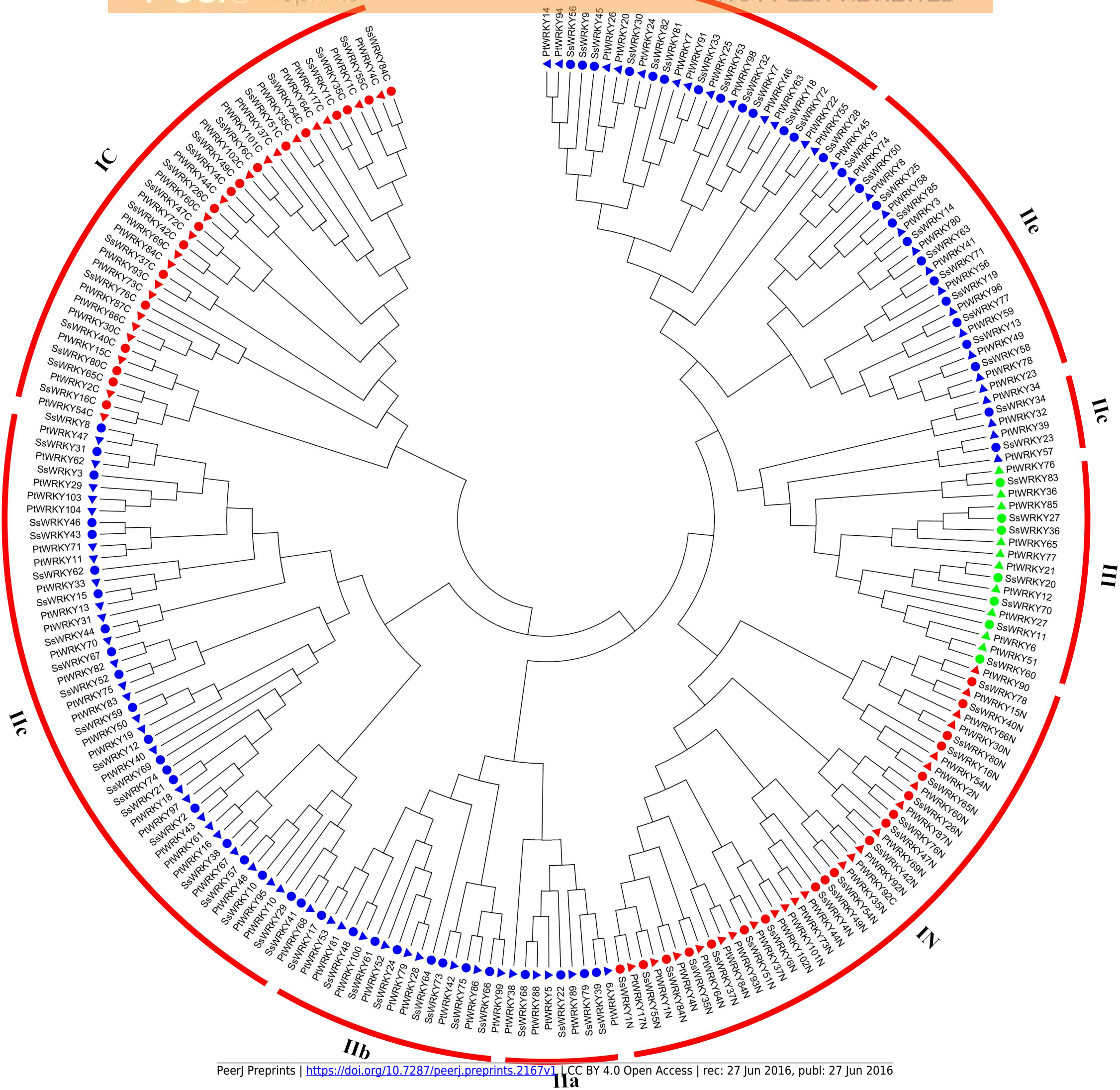


## Figure 4 (on next page)

### Phylogenetic tree of WRKY domains from willow and poplar.

The phylogenetic tree was constructed using the neighbor-joining method in MEGA 6.0. The WRKY genes with the suffix 'N' and 'C' indicate the N-terminal and the C-terminal WRKY domains of group I, respectively. The different colors indicate different groups (I, II and III) or subgroups (IIa, b, c, d and e) of WRKY domains. Circles indicate WRKY genes from willow, and triangles represent genes from poplar.

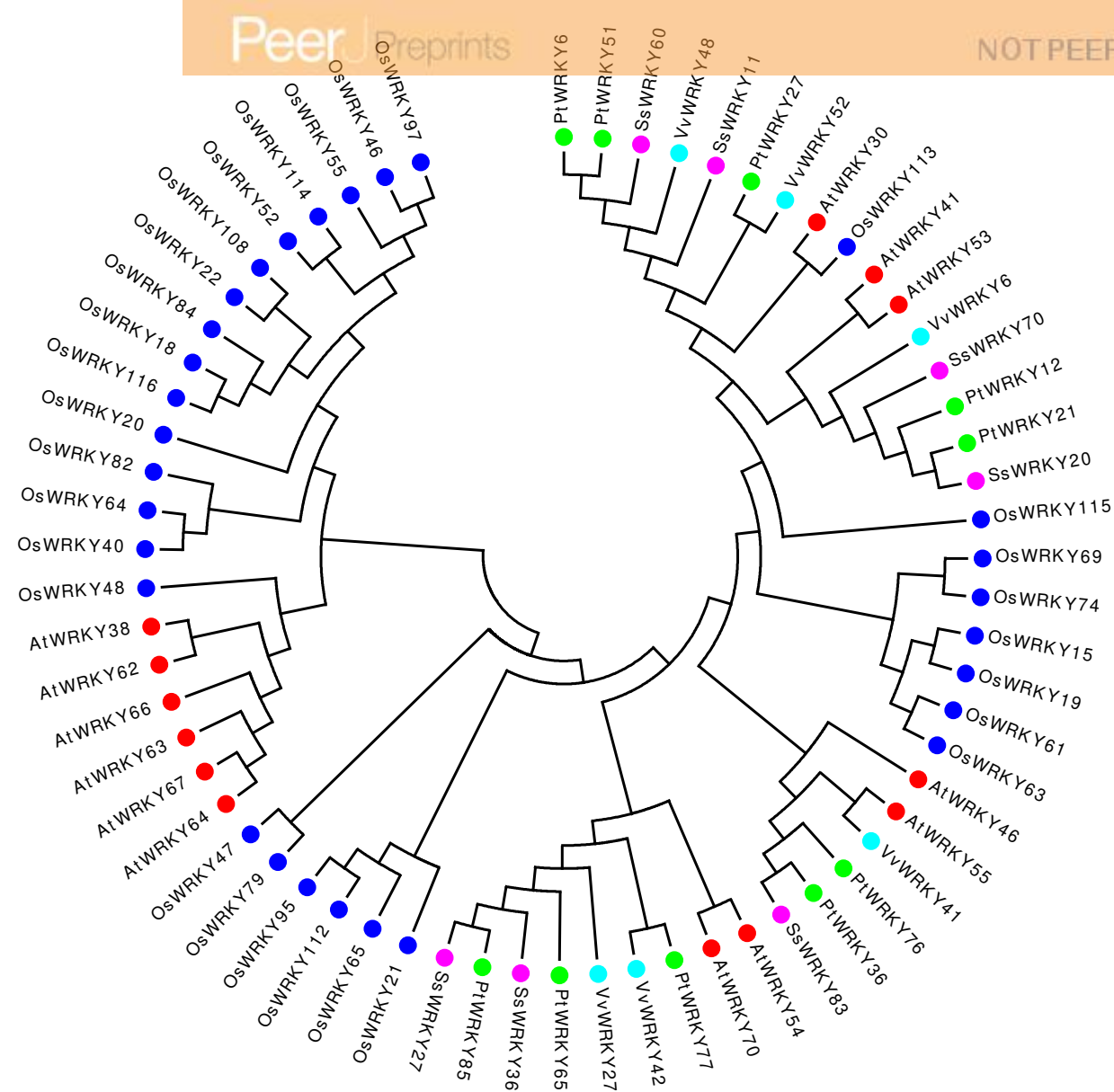




## Figure 5 (on next page)

**Phylogenetic tree of full-length group III WRKY genes from *Arabidopsis* (AtWRKY), rice (OsWRKY), grape (VvWRKY), poplar (PtWRKY) and willow (SsWRKY).**

The phylogenetic tree was constructed using the neighbor-joining method in MEGA 6.0. Dicotyledonous (*Arabidopsis*, grape, poplar and willow) and monocotyledonous (rice) WRKY III genes are marked with colored dots.

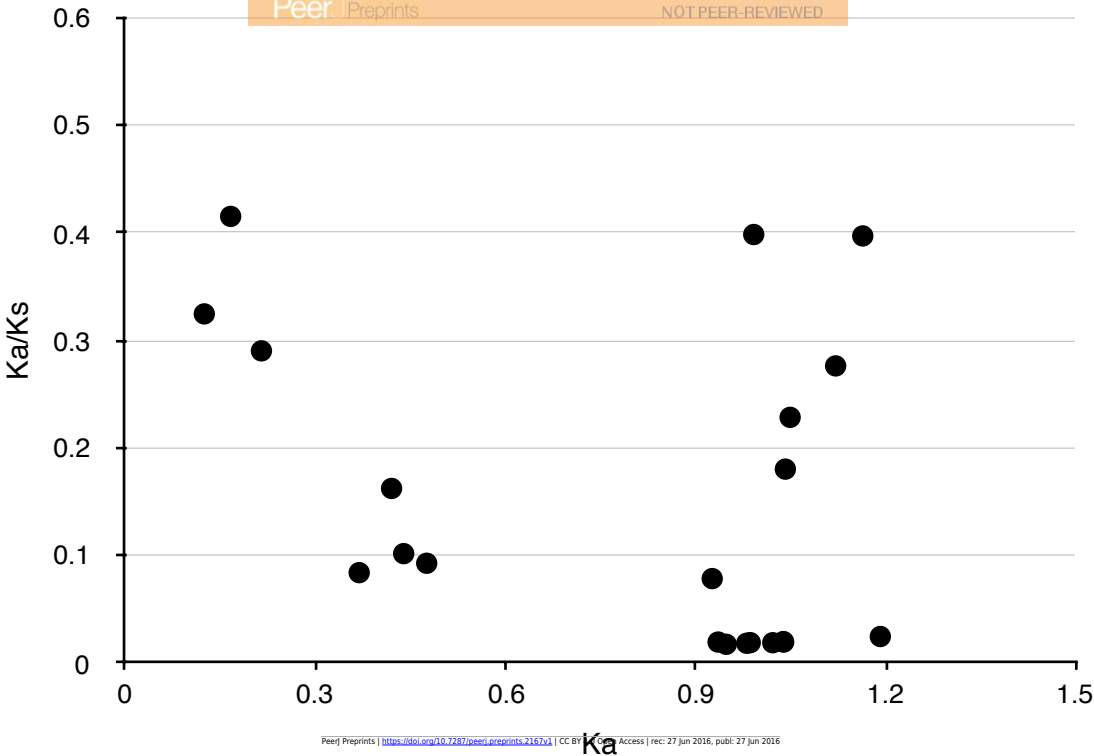




# Figure 6(on next page)

## Scatter plots of the Ka/Ks ratios of WRKY III genes in willow.

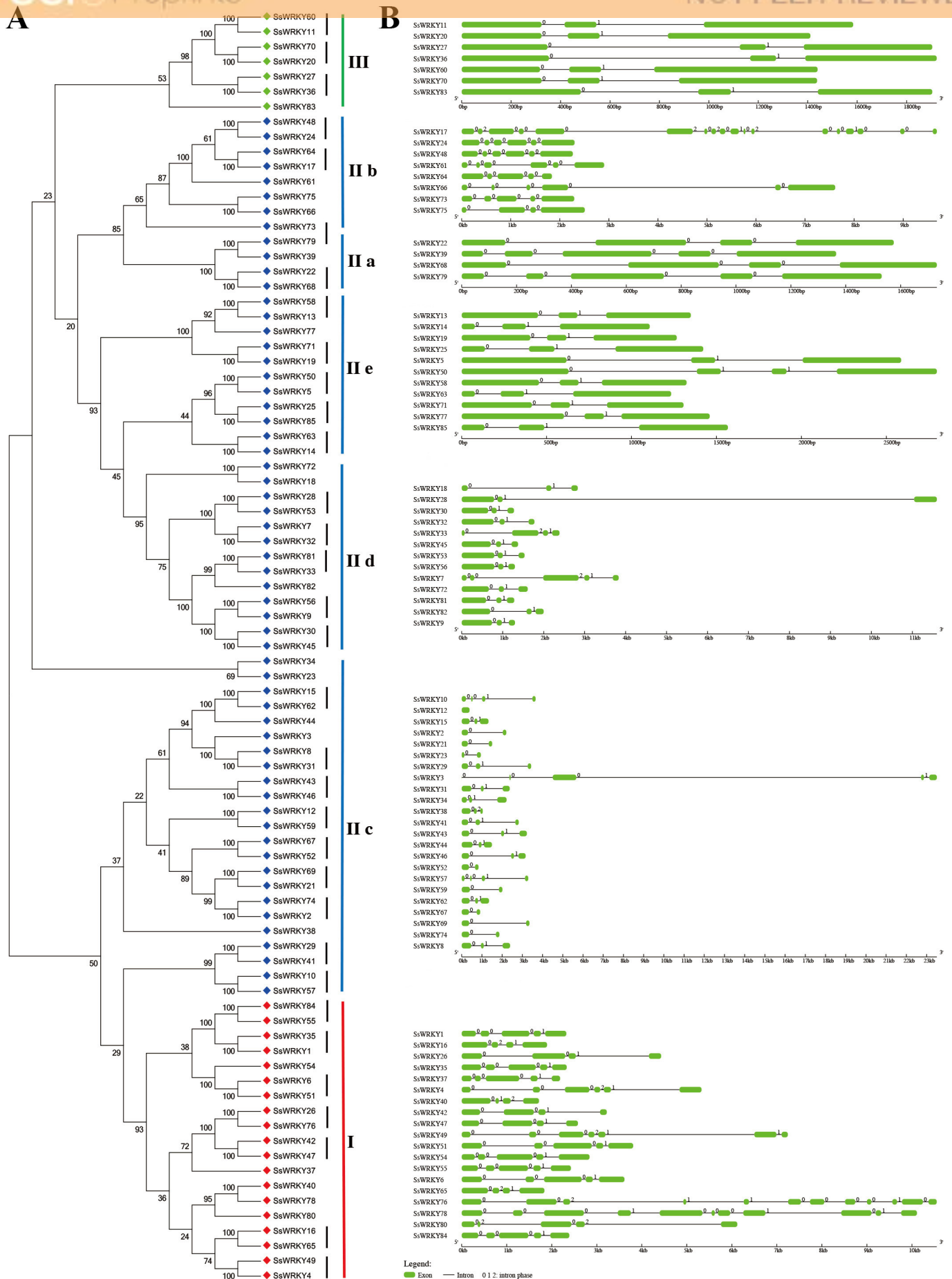
The Y- and X-axes denote the Ka/Ks ratio and Ka for each pair, respectively.



## Figure 7 (on next page)

### Genomic organization of SsWRKY genes.

(A) The phylogenetic tree built on the basis of full-length SsWRKY genes was depicted using the neighbor-joining method in MEGA 6.0. The short black lines indicate the existence of duplicated gene pairs; (B) The graphic exon-intron structure of SsWRKY genes is displayed using GSDS. Green indicates exons, and gray indicates introns. The introns phases 0, 1 and 2 are indicated by numbers 0, 1 and 2, respectively.



## Figure 8(on next page)

**The distribution of twenty conserved motifs of SsWRKY genes was identified by the online program MEME.**

The names of all members are displayed on the left side of the figure. Different motifs are displayed in different colored boxes as indicated on the right side. The conserved motifs 1, 2, 3, and 5, broadly distributed across SsWRKY genes, were definitely characterized as the WRKY conserved domains.

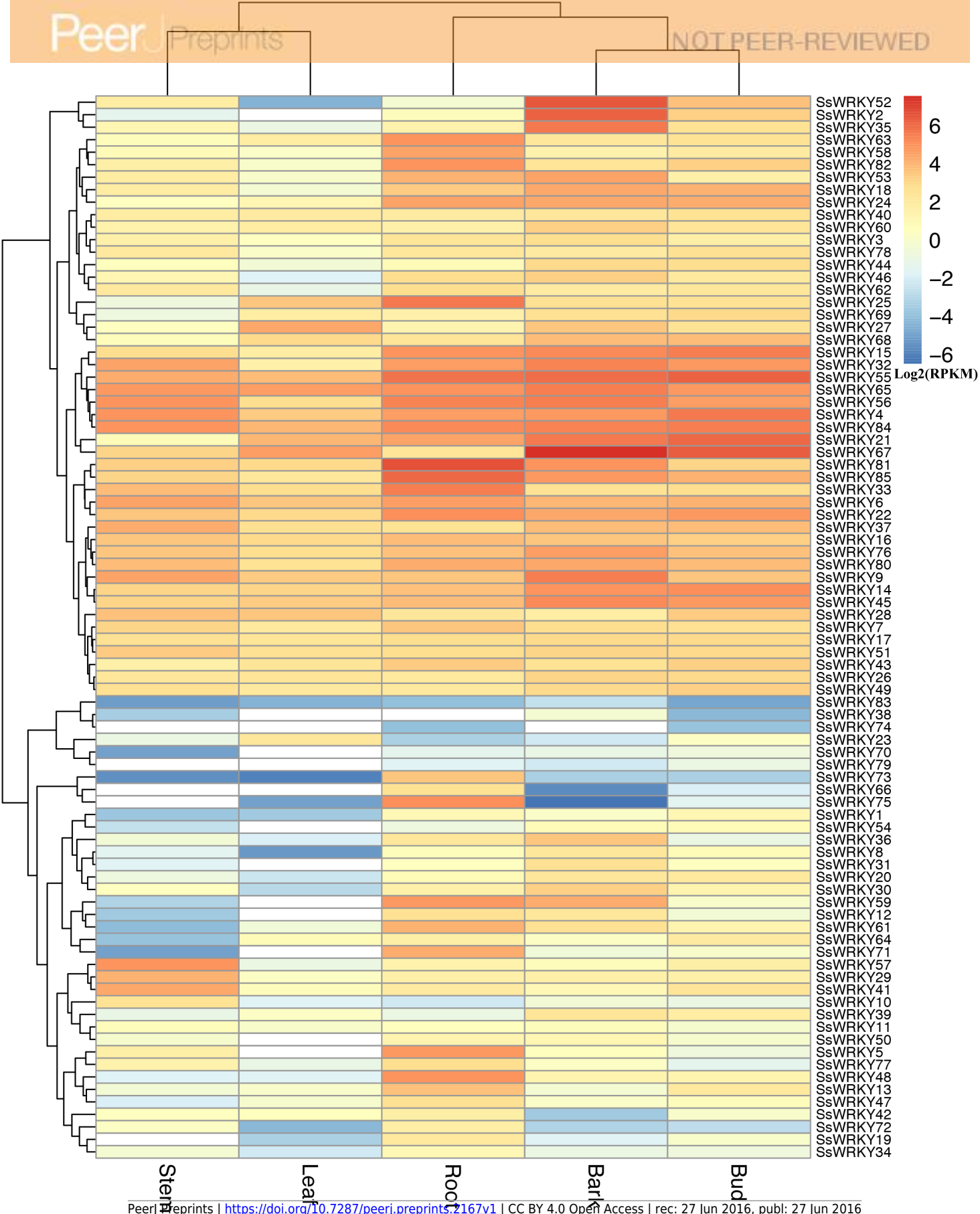


## Figure 9 (on next page)

### Expression profiles of the 85 SsWRKY genes in root, stem, bark, bud and leaf.

Color scale represents RPKM normalized log2 transformed counts and red indicates high expression, blue indicates low expression and white indicates the gene is not expressed in this tissue.





# **Table 1**(on next page)

**The detailed characteristics of WRKY genes identified in willow.**

Gene	SequenceID	Chr	Group	Ortholog		Deduced polypeptide			Introns
				AtWRKY	PtWRKY	Length(aa)	PI	MW(kDa)	
SsWRKY1	willow_GLEAN_10011238	1	I	33	17	583	7.14	64.7	4
SsWRKY2	willow_GLEAN_10019192	1	II c	45	43	162	9.47	18.6	1
SsWRKY3	willow_GLEAN_10017208	1	II c	28,71	29	584	9.42	65.6	4
SsWRKY4	willow_GLEAN_10017139	1	I	20	44	560	6.99	60.9	5
SsWRKY5	willow_GLEAN_10007860	1	II e	35	45	445	5.92	48.4	2
SsWRKY6	willow_GLEAN_10003806	1	I	2	37,101,102	733	5.69	78.8	4
SsWRKY7	willow_GLEAN_10022392	2	II d	21	46,63	453	9.53	49.9	4
SsWRKY8	willow_GLEAN_10022273	2	II c	71	47	328	6.89	37.0	2
SsWRKY9	willow_GLEAN_10009329	2	II d	15	14,94	339	9.77	37.5	2
SsWRKY10	willow_GLEAN_10009231	2	II c	12	48	204	7.64	23.6	3
SsWRKY11	willow_GLEAN_10016913	2	III	30	6,51	351	6.27	39.2	2
SsWRKY12	willow_GLEAN_10016886	2	II c	-	19,50	129	6.75	14.6	0
SsWRKY13	willow_GLEAN_10016883	2	II e	22	23,49,78	352	5.81	38.3	2
SsWRKY14	willow_GLEAN_10019911	2	II e	-	3	247	5.58	28.1	2
SsWRKY15	willow_GLEAN_10019925	2	II c	23	13,33	319	6.46	35.6	2
SsWRKY16	willow_GLEAN_10019982	2	I	1	54	472	6.88	52.2	3
SsWRKY17	willow_GLEAN_10020022	2	II b	47	53	1081	5.25	116.8	17
SsWRKY18	willow_GLEAN_10025583	3	II d	-	55	142	9.60	16.5	2
SsWRKY19	willow_GLEAN_10025423	3	II e	29	41	335	5.54	37.9	2
SsWRKY20	willow_GLEAN_10025378	3	III	41/53	21	342	5.25	38.4	2
SsWRKY21	willow_GLEAN_10008020	3	II c	45	18	157	9.41	17.8	1
SsWRKY22	willow_GLEAN_10006448	3	II a	40	88	320	8.38	35.4	3
SsWRKY23	willow_GLEAN_10013342	3	II c	-	39	109	8.03	12.9	1
SsWRKY24	willow_GLEAN_10009960	4	II b	42	28,79	604	6.93	65.3	5
SsWRKY25	willow_GLEAN_10017267	4	II e	65	8,58	267	5.43	29.7	2
SsWRKY26	willow_GLEAN_10018559	4	I	58	60	537	8.72	58.9	3
SsWRKY27	willow_GLEAN_10004854	4	III	54	85	323	5.70	36.3	2
SsWRKY28	willow_GLEAN_10008312	5	II d	-	-	490	10.27	54.0	2
SsWRKY29	willow_GLEAN_10009112	5	II c	13	68	235	8.70	26.7	2
SsWRKY30	willow_GLEAN_10003565	5	II d	15	20	310	9.48	34.3	2
SsWRKY31	willow_GLEAN_10016009	5	II c	28,71	62	322	6.67	36.2	2
SsWRKY32	willow_GLEAN_10018195	5	II d	21	46,63	349	9.69	38.8	2
SsWRKY33	willow_GLEAN_10026833	6	II d	7	91	339	9.89	36.8	3
SsWRKY34	willow_GLEAN_10026721	6	II c	49	34	287	5.25	32.1	2
SsWRKY35	willow_GLEAN_10026591	6	I	33	64	572	6.41	62.7	4
SsWRKY36	willow_GLEAN_10026566	6	III	54	85	329	6.13	36.7	2
SsWRKY37	willow_GLEAN_10020588	6	I	44	93	478	9.25	52.5	4
SsWRKY38	willow_GLEAN_10026166	6	II c	51	67	233	5.03	26.1	2

SsWRKY39	willow_GLEAN_10026455	6	II a	18/60	9	327	9.02	36.2	4
SsWRKY40	willow_GLEAN_10026458	6	I	32	15	413	8.26	44.9	3
SsWRKY41	willow_GLEAN_10008192	7	II c	13	68	236	9.21	26.6	2
SsWRKY42	willow_GLEAN_10025108	8	I	3/4	69	460	8.80	50.6	3
SsWRKY43	willow_GLEAN_10025123	8	II c	57	71	295	6.32	32.3	2
SsWRKY44	willow_GLEAN_10015641	8	II c	48	70	357	6.11	39.9	2
SsWRKY45	willow_GLEAN_10008155	9	II d	15	20,26	331	9.57	36.4	2
SsWRKY46	willow_GLEAN_10013562	10	II c	57	71	289	6.26	31.9	2
SsWRKY47	willow_GLEAN_10013586	10	I	3/4	72	490	8.60	53.7	3
SsWRKY48	willow_GLEAN_10004012	11	II b	42	100	585	6.48	63.3	5
SsWRKY49	willow_GLEAN_10006060	11	I	20	44	607	7.09	6.6	6
SsWRKY50	willow_GLEAN_10007614	11	II e	35	74	481	5.39	51.6	3
SsWRKY51	willow_GLEAN_10007542	11	I	2	37	734	6.10	79.7	4
SsWRKY52	willow_GLEAN_10013801	12	II c	-	75	178	9.08	20.5	1
SsWRKY53	willow_GLEAN_10012158	13	II d	74	25	356	9.66	40.0	2
SsWRKY54	willow_GLEAN_10004417	13	I	2	35	697	6.52	76.1	4
SsWRKY55	willow_GLEAN_10007732	13	I	33	1	602	7.65	66.0	4
SsWRKY56	willow_GLEAN_10009039	14	II d	15	14,94	362	9.39	40.0	2
SsWRKY57	willow_GLEAN_10016668	14	II c	12	48	180	8.47	20.7	3
SsWRKY58	willow_GLEAN_10016177	14	II e	22	23,49,78	354	6.35	38.8	2
SsWRKY59	willow_GLEAN_10016180	14	II c	43	19,50	193	9.47	21.7	1
SsWRKY60	willow_GLEAN_10016220	14	III	30	6	368	5.03	41.3	2
SsWRKY61	willow_GLEAN_10018940	14	II b	42	28,79	467	8.78	50.0	5
SsWRKY62	willow_GLEAN_10018891	14	II c	23	13,33	318	5.71	35.6	2
SsWRKY63	willow_GLEAN_10018881	14	II e	-	80	263	5.05	29.7	2
SsWRKY64	willow_GLEAN_10020302	14	II b	36	-	460	6.28	50.0	4
SsWRKY65	willow_GLEAN_10020380	14	I	1	2	481	5.98	52.8	3
SsWRKY66	willow_GLEAN_10011119	15	II b	9	99	618	6.55	66.2	5
SsWRKY67	willow_GLEAN_10016438	15	II c	-	82	178	9.35	20.5	1
SsWRKY68	willow_GLEAN_10023347	16	II a	40	88	320	8.82	35.3	3
SsWRKY69	willow_GLEAN_10023447	16	II c	45	18	178	9.17	20.1	1
SsWRKY70	willow_GLEAN_10023687	16	III	41/53	21	336	5.17	37.2	2
SsWRKY71	willow_GLEAN_10023735	16	II e	29	41	325	5.54	36.6	2
SsWRKY72	willow_GLEAN_10014752	16	II d	-	55	338	9.24	37.9	2
SsWRKY73	willow_GLEAN_10009602	16	II b	9	42	509	5.51	55.3	4
SsWRKY74	willow_GLEAN_10010473	17	II c	45	43	182	9.92	20.9	1
SsWRKY75	willow_GLEAN_10015128	17	II b	9	86	544	6.01	59.0	3
SsWRKY76	willow_GLEAN_10015184	17	I	58	87	1044	8.94	116.1	11
SsWRKY77	willow_GLEAN_10005468	17	II e	27	96	411	5.96	45.7	2
SsWRKY78	willow_GLEAN_10006860	18	I	-	90	1593	8.67	179.0	10

SsWRKY79	willow_GLEAN_10006862	18	II a	18/60	9	320	8.57	35.6	4
SsWRKY80	willow_GLEAN_10011608	18	I	32	-	528	5.74	57.8	4
SsWRKY81	willow_GLEAN_10004546	18	II d	7	7,91	300	9.80	32.8	2
SsWRKY82	willow_GLEAN_10003422	19	II d	11/17	24	339	9.58	37.1	2
SsWRKY83	willow_GLEAN_10011321	19	III	55	36,76	358	5.63	38.7	2
SsWRKY84	willow_GLEAN_10005288	19	I	33	4	597	6.69	65.6	4
SsWRKY85	willow_GLEAN_10002834	N/A	II e	65	58	268	5.83	30.2	2

Chr, chromosome numbers.

N/A, not available.

"-", not detected.

## Table 2 (on next page)

The details of twenty conserved motif sequences identified in SsWRKY genes.

Motif	Width	Best possible match
1	29	ILDDGYRWRKYGQKVIKGNPYPRSYRCT
2	29	CPVRKHVERCWEDPTMVITTYEGEHNHPW
3	37	PSDDGYNWRKYGQKQVKGSEYPRSYKCTHPNCPVKK
4	21	KKGHKKIREPRFAFQTRSEVD
5	29	KVECSHDGHITEIYKGTNHHPKPQPNCR
6	15	KRRKNRVKWVVRVPA
7	50	KEELAVLQEELNRMKEENKRLKEMLDQICENYNALQMHFMDLMQQNNE
8	29	PVIRSPYFTIPPGLSPTTELDSPVFFSNS
9	29	LVEQMTAAITADPNFTAALAAASGIMGQ
10	28	QVQYRNCMVITDETVFKFKKVISLLNRT
11	29	LQQQQQQQMKYQADMMYRKSNSGINLNF
12	15	MRKARVSVRARCEAP
13	50	MDGTVANLDGDAFHLMGMPHSSDHISQQHKRKCGRGEDGNVKGCGSSGI
14	21	PPAAMAMASTTSAAASMLLSG
15	21	VEEAARAGIESCEHVIRLLCQ
16	21	MATISASAPFPTITLDTQNP
17	40	LGHGRVRKLKLP SHLPQNIFLDNPHCKTIHAPKPPQMVP
18	17	LLPDYGLLQDIVPSMH
19	17	GGEDDEDEPEPKRWKIE
20	49	PSPTTGTFPGQAFNWKSNSGDNQQGVKGEDKDFSDFSFQTPARPPATSS



# Table 3 (on next page)

The number of WRKY genes identified in *Arabidopsis thaliana*, *Cucumis sativus*, *Populus trichocarpa*, *Vitis vinifera*, *Salix suchowensis* and *Oryza sativa*.

Species	Group						
	I	IIa	IIb	IIc	IId	IIe	III
<i>Arabidopsis thaliana</i>	13	4	7	18	7	9	14
<i>Cucumis sativus</i>	10	4	4	16	8	7	6
<i>Populus trichocarpa</i>	50	5	9	13	13	4	10
<i>Vitis vinifera</i>	12	4	8	16	7	6	6
<i>Salix suchowensis</i>	19	4	8	23	13	11	7
<i>Oryza sativa</i>	34	4	8	7	11	0	36