

# The WRKY superfamily of plant transcription factors

Thomas Eulgem, Paul J. Rushton, Silke Robatzek and Imre E. Somssich

**The WRKY proteins are a superfamily of transcription factors with up to 100 representatives in *Arabidopsis*. Family members appear to be involved in the regulation of various physiological programs that are unique to plants, including pathogen defense, senescence and trichome development. In spite of the strong conservation of their DNA-binding domain, the overall structures of WRKY proteins are highly divergent and can be categorized into distinct groups, which might reflect their different functions.**

One of the apparent fundamental principles of biological evolution is that the progression from ancient to advanced life forms is inseparably connected to an increase in regulatory capacity. Genome-sequencing efforts have provided evidence for a positive correlation between the proportion of genes involved in information processing and the complexity of organisms. More than 20% of the genes within the sequence available for the *Arabidopsis thaliana* genome appear to encode proteins that play a role in signal transduction or transcription<sup>1</sup>, whereas only 12% of the genome of the single-celled yeast *Saccharomyces cerevisiae* contains genes of this type<sup>2</sup>.

This increase in biological complexity coincides with the appearance or expansion of specific groups of regulator genes. One example is the nuclear-receptor-gene family, which is completely absent in yeast but highly represented in metazoan organisms<sup>3</sup>. The evolution of nuclear receptors is believed to be a key event in the development of intercellular communication, a prerequisite for the multicellularity of metazoans<sup>4</sup>. Similarly, the establishment of a complex animal body plan was driven by the amplification and divergence of ancestral homeobox genes, thereby generating a sophisticated regulatory system of functionally interconnected transcriptional regulators<sup>5</sup>.

To meet their disparate biological requirements, plants and animals have evolved unique regulatory mechanisms. This was partly achieved by combining functional domains from pre-existing factors to build new regulators, as exemplified by the MADS-box factors, which play a central role in determining floral and organ identity in plants<sup>6</sup>. In addition, completely new factors have arisen and we focus here on the potential biological roles of WRKY (pronounced 'worky') proteins, a large family of transcriptional regulators that has to date only been found in plants. The abundance of information provided by the *Arabidopsis* sequencing projects is an ideal basis for comparative analysis of this superfamily within one plant species. Although their precise regulatory functions are largely unknown, the fact that these factors appear to be specific to plants, with probably up to 100 members in *Arabidopsis*, suggests that they play an important role during plant evolution.

## Biochemical properties of WRKY proteins

The first WRKY cDNAs were cloned from sweet potato (*Ipomoea batatas*; SPF1), wild oat (*Avena fatua*; ABF1,2), parsley (*Petroselinum crispum*; PcWRKY1,2,3) and *Arabidopsis* (ZAP1), based on the ability of the encoded proteins to bind specifically to the DNA sequence motif (T)(T)TGAC(C/T), which is known as the W box<sup>7-10</sup>. It has been suggested that the cognate binding site for SPF1 is different from other WRKY proteins. However, the oligonucleotide used to isolate SPF1 does have a W box in the flanking sequence<sup>7</sup>.

The name of the WRKY family is derived from the most prominent feature of these proteins, the WRKY domain, a 60 amino acid region that is highly conserved amongst family members. The emerging picture is that these proteins are regulatory transcription factors with a binding preference for the W box, but with the potential to differentially regulate the expression of a variety of target genes. Consistent with a role as transcription factors, PcWRKY1 and WIZZ (from tobacco) have been shown to be targeted to the nucleus<sup>11,12</sup>.

## The WRKY domain and the W box

The WRKY domain is defined by the conserved amino acid sequence WRKYGQK at its N-terminal end, together with a novel zinc-finger-like motif<sup>8</sup> (Fig. 1). Because of the clear binding preference of all characterized WRKY proteins for the same DNA motif, it has been assumed that the WRKY domain, as their only conserved structural feature, constitutes a DNA-binding domain. Indeed, it has recently been shown that an isolated WRKY domain has sequence-specific DNA-binding activity<sup>12</sup>. The divalent metal chelators 1,10-*o*-phenanthroline and EDTA abolish *in vitro* DNA binding, which is taken as strong support for a zinc-finger structure within the WRKY domain<sup>8,10,11</sup>. However, it has not yet been proven that zinc is actually complexed in the WRKY domain. In addition, nothing is known about the function of the WRKYGQK heptapeptide stretch, the hallmark of this superfamily.

All known WRKY proteins contain either one or two WRKY domains. They can be classified on the basis of both the number of WRKY domains and the features of their zinc-finger-like motif. WRKY proteins with two WRKY domains belong to group I, whereas most proteins with one WRKY domain belong to group II (Fig. 2). Generally, the WRKY domains of group I and group II members have the same type of finger motif, whose pattern of potential zinc ligands (C-X<sub>4-5</sub>-C-X<sub>22-23</sub>-H-X<sub>1</sub>-H; Fig. 1) is unique among all described zinc-finger-like motifs<sup>13</sup>. The single finger motif of a small subset of WRKY proteins is distinct from that of group I and II members. Instead of a C<sub>2</sub>-H<sub>2</sub> pattern, their WRKY domains contain a C<sub>2</sub>-HC motif (C-X<sub>7</sub>-C-X<sub>23</sub>-H-X<sub>1</sub>-C; Fig. 1). Owing to this distinction, they were recently assigned to the newly defined group III. Nevertheless, experimental evidence has shown that members of all three groups bind sequence specifically to various W box elements (R.S. Cormack *et al.*, unpublished).

The two WRKY domains of group I members appear to be functionally distinct. As has been shown for SPF1, ZAP1 and PcWRKY1, sequence-specific binding to their target DNA sequences is mediated mainly by the C-terminal WRKY domain<sup>7,10,12</sup>. The function of the N-terminal WRKY domain remains unclear. Because protein regions outside of the C-terminal WRKY domain contribute to the overall strength of DNA

Group I

WRKY1 TLFDIVNDGYRWRKYGQKSVKGSYPYRYSYRCSSPG...CPVKKHVERSSHDTKLLITTYEGKEDHDM  
 WRKY2 SDVDLDDGYRWRKYGQKVVKGNPNPRYSYKCTAPG...CTVRKHVERASHDLKSVITTYEGKENDVP  
 WRKY3 SEVDLDDGYRWRKYGQKVVKGNPNPRYSYKCTTPD...CGVRKHVERAATDPKAVTTYEGKENDVP  
 WRKY4 SEVDLDDGYRWRKYGQKVVKGNPNPRYSYKCTTPG...CGVRKHVERAATDPKAVTTYEGKENDVP  
 WRKY20 SEVDLDDGYRWRKYGQKVVKGNPNPRYSYKCTAHG...CPVRKHVERASHDPKAVITTYEGKEDHDM  
 WRKY25 SDIDLVDLIDGFRWRKYGQKVVKGNPNPRYSYKCTFQG...CGVKKQVERSAADERAVLTTYEGRENHDI  
 WRKY26 SDIDLDDGYRWRKYGQKVVKGNPNPRYSYKCTFTG...CFVRKHVERAFQDPKSVITTYEGKEDHDM  
 WRKY32 GDVIGCGDGYRWRKYGQKVVKGNPNPRYSYKCTSAG...CPVRKHVERASHDRAVITTYEGKENDVP  
 WRKY33 SDIDLDDGYRWRKYGQKVVKGNPNPRYSYKCTTIG...CPVRKHVERASHDMRAVITTYEGKENDVP  
 WRKY34 SDIDLDDGYRWRKYGQKVVKGNPNPRYSYKCTANG...CTVTKHVERASDDFKSVLTTYIGKEDHDM  
 WRKY44 VESDSLEDGFRWRKYGQKVVGNAYPRYSYKCTSAN...CRARKHVERASDDPRAFITTYEGKENDHLL  
 WRKY45 SQVDLDDGYRWRKYGQKAVKNNPFPRYSYKCTEEG...CRVKKQVQRQWGDQVVTTYQGVHTHVD  
 WRKY58 SEVDLDDGYRWRKYGQKVVKGNPNPRYSYKCTTPN...CTVRKHVERASTDAKAVITTYEGKENDVP  
 WRKY10 SDEDNPNDDGYRWRKYGQKVVKGNPNPRYSYKCTNIE...CRVKKHVERGADNIKLVTTYDGIENDHPS

Group II

(a) WRKY18 DTSLTVKDGFRWRKYGQKVVTRDNPSPRAYFRCSFAPS...CPVKKKQVRSADPSSLVATYEGTENDHLP  
 WRKY40 KDGFRWRKYGQKVVTRDNPSPRAYFRCSFAPS...CSVKKKQVRSVEDQSVLVA TYEGENDHMP  
 WRKY60 VSSLTVKDGFRWRKYGQKVVTRDNPSPRAYFRCSFAPS...CLVKKKQVRSADPSSLVATYEGTENDHLP

(b) WRKY6 SEAPMISDGCQWRKYGQKMAKGNPCPRAYYRCMTATG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY9 CETATMNDGCQWRKYGQKMAKGNPCPRAYYRCMTVAPG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY31 SEAAAMISDGCQWRKYGQKMAKGNPCPRAYYRCMTMAGG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY36 CEDPISDGCQWRKYGQKMAKGNPCPRAYYRCMTMAGG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY42 SEAPMISDGCQWRKYGQKMAKGNPCPRAYYRCMTAVG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY47 HKQHEVNDGCQWRKYGQKMAKGNPCPRAYYRCMTAVG...CPVRKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY61 NDGCQWRKYGQKMAKGNPCPRAYYRCMTIAAS...CPVRKQVQRCAEDRSILITTYEGNENDHPL

(c) WRKY8 TEVDHLEDGYRWRKYGQKAVKNSPYRYSYRCCTTQK...CNVKKRVERSYQDPTVVITTYEQENDHPI  
 WRKY12 SDVDVLDGFRWRKYGQKVVKNSLHPRYSYRCCTHNN...CRVKKRVERLSEDCRMVITTYEGRENHIPS  
 WRKY13 SEVDVLDGFRWRKYGQKAVKNSPYRYSYRCCTQDK...CRVKKRVERLADPRMVITTYEGRENHIPS  
 WRKY23 SEVDHLEDGYRWRKYGQKAVKNSPYRYSYRCCTTAS...CNVKKRVERSFQDPTVVITTYEQENDHPI  
 WRKY24 SDDVDVLDGFRWRKYGQKAVKNSPYRYSYRCCTYHT...CNVKKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY28 SEVDHLEDGYRWRKYGQKAVKNSPYRYSYRCCTTQK...CNVKKRVERSFQDPTVVITTYEQENDHPI  
 WRKY43 SDADILDGFRWRKYGQKAVKNSPYRYSYRCCTQHM...CNVKKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY48 KSIDNLDGFRWRKYGQKAVKNSPYRYSYRCCTTVG...CGVKKRVERSSDPSIVITTYEQENDHPI  
 WRKY49 NSNGMCDGFRWRKYGQKAVKNSPYRYSYRCCTNPI...CNAKKQVERSIDESNTYIITTYEGFENDHPI  
 WRKY50 SEVEVLDGFRWRKYGQKAVKNSPYRYSYRCCTVHG...CPVKKRVERDRDPSFVITTYEGSENHESM  
 WRKY51 DVMDGFRWRKYGQKAVKNSPYRYSYRCCTSEG...CSVKKRVERDGDAAAYVITTYEGVENHESL  
 WRKY56 SDDVDVLDGFRWRKYGQKAVKNSPYRYSYRCCTYHT...CNVKKQVQRCAEDRSILITTYEGNENDHPL  
 WRKY57 SDVDNLEDGFRWRKYGQKAVKNSPYRYSYRCCTNSR...CTVKKRVERSSDPSIVITTYEQENDHPI  
 WRKY59 DEKVALDDGFRWRKYGQKAVKNSPYRYSYRCCTSPD...CNVKKKIERDTNPNPDYILT TYEGRENHIPS

(d) WRKY7 KMADIPDDEFWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERALDDAMMLIVTYEGENDHALV  
 WRKY11 KIADIPDDEFWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERALDDPAMLIVTYEGENDHNS  
 WRKY15 KMSDVPDDYSWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERAADSSMLIVTYEGENDHNS  
 WRKY17 KIADIPDDEFWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERALDDSTMLIVTYEGENDHNS  
 WRKY21 KVADIPDDEFWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERALDDPAMLIVTYEGENDHNS  
 WRKY39 KIADIPDDEFWRKYGQKPIKGSPPHPRGYKCSVVRG...CPARKHVERCIDETSMLIVTYEGENDHNS

(e) WRKY14 SGEVVPDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CSARKQVERSRTPDNMLVITYTSENHNPWP  
 WRKY16 DRGSRSSDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CFARKQVERSRTPDNVSVITYISENHPFP  
 WRKY22 AAALNSDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CLARKQVERNRSDPKMFIVITYTAEENHPAP  
 WRKY27 TQENLSSDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CLARKQVERNLDPNIFIVITYTGEENHPRP  
 WRKY29 KEENLSSDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CLARKQVERNPNQPEKFTITYTNEENHELP  
 WRKY35 SGEVVPDLWAWRKYGQKPIKGSPPHPRGYKCSVVRG...CSARKQVERSRTPDNMLVITYTSENHNPWP

Group III

WRKY30 GVDRTILDGFRWRKYGQKILGAKFPRGYKCTYRKSQG...EATKQVQRSDENQMLLEISYRGIENDHNSQA  
 WRKY41 GLEGPHDDIFSWRKYGQKILGAKFPRYSYKCTFRNTQY...CWATKQVQRSDGPTIFEVTYRGTENDHNSQ  
 WRKY46 QENGSIDGHCWRKYGQKEIHGSKNPRAYYRCCTHRFTQD...CLAVKQVQRSDTDPPLFVYKLGNDHNSNI  
 WRKY53 GLEGPQDDVFSWRKYGQKILGAKFPRYSYKCTHRSTQD...CWATKQVQRSDGPTIFEVTYRGTENDHNSQA  
 WRKY54 VEAKSSEDRYAWRKYGQKEILNTTFPRYSYKCTHKTQ...CKATKQVQRSDSEMFIITYIGYHTENDHNS  
 WRKY55 NTDLPDNDHNTWRKYGQKEILGSRFPAYYRCCTHQLY...CPAKKQVQRSDPFTFRVYRGTENDHNS

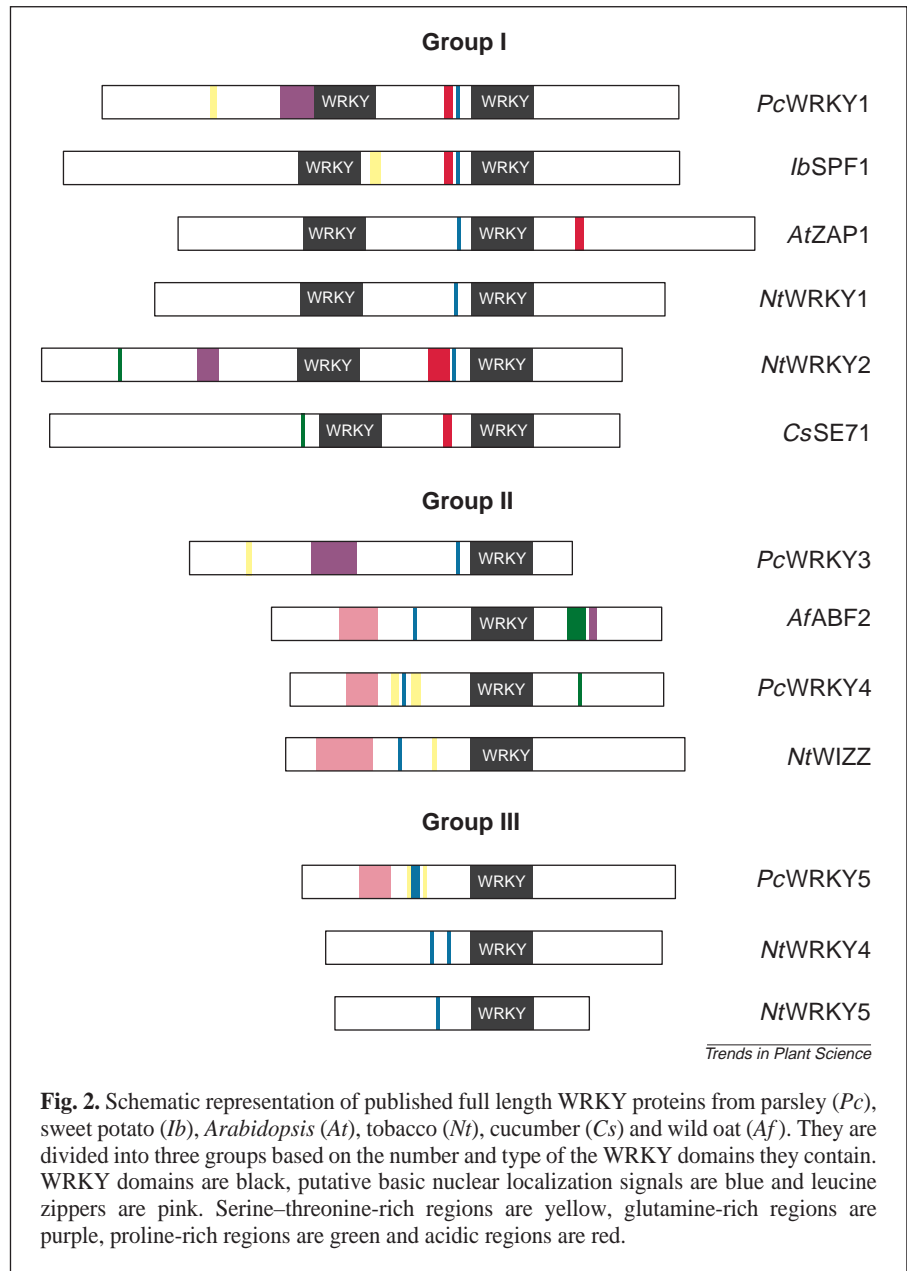
WRKY38 SPDPIYYDGYLWRKYGQKIKKSNHQRYSYKCSYNKDN...CEARKHEQIKDNPPVYRTTYFGHTENDHNSKTE  
 WRKY52 IPAIDEGDLWTWRKYGQKILGSRFPAYYRCCTHKTQ...CKATKQVQRSDSEMFIITYIGYHTENDHNS

**Fig. 1.** Left Comparison of WRKY domain sequences from *At*WRKY proteins. Sequences encoding the peptide stretch WRKYGQK were found by the BLAST programs tblastn and blastp programs<sup>37</sup> in genomic and EST databases. Gaps (dots) have been inserted for optimal alignment. Residues that are highly conserved within each of the major groups are in red and potential zinc ligands are highlighted in black boxes. For each (sub)group, the position of a conserved intron is indicated by an arrowhead.

binding, the N-terminal domain might participate in the binding process, increasing the affinity or specificity of these proteins for their target sites. Alternatively, it might provide an interface for protein–protein interactions, a known function of some zinc-finger-like domains<sup>14</sup>; this could allow more efficient DNA binding through interactions with other DNA-associated proteins. Not unexpectedly, the single WRKY domains of group II and III family members are more similar in sequence to the C-terminal than to the N-terminal WRKY domain of group I proteins, suggesting that the C-terminal and single WRKY domains are functionally equivalent and constitute the major DNA-binding domain.

The conservation of the WRKY domain is mirrored by a remarkable conservation of the cognate *cis*-acting W box elements. These (T)(T)TGAC(C/T) sequence elements contain the invariant TGAC core, which is essential for function and WRKY binding. They mediate transcriptional responses to pathogen-derived elicitors<sup>9,15</sup> and are present in the promoters of many plant genes that are associated with defense<sup>16</sup>. Functional W boxes frequently cluster within short promoter stretches<sup>15–17</sup> and can act together synergistically<sup>12</sup>. WRKY–W box interactions have been demonstrated by numerous binding experiments, both *in vitro* and *in vivo*<sup>8–10,12,18,19</sup> (R.S. Cormack *et al.*, unpublished), and random binding-site selection assays have shown that the optimal binding site for ZAP1 contains the W box motif<sup>10</sup>. Interactions of WRKY proteins with W boxes can be regulated post-translationally, because binding of WRKY-like DNA-binding activities to W boxes in tobacco is abolished by treatment with alkaline phosphatase<sup>18</sup> and the protein-kinase inhibitor staurosporin<sup>20</sup>.

In spite of the stereotypic binding preferences of WRKY proteins for W boxes, their affinities for certain types or arrangements of this element can vary (R.S. Cormack *et al.*, unpublished). Sequences flanking the invariant W box TGAC core might be partly responsible for the observed specificity. In addition, the cooperative assembly of discrete higher-order WRKY–DNA complexes at defined W box arrangements might also account for specific promoter recognition<sup>12</sup>. Owing to the high variability in overall protein structure, access to certain promoters would be restricted to distinct family members that fit into the three-



**Fig. 2.** Schematic representation of published full length WRKY proteins from parsley (*Pc*), sweet potato (*Ib*), *Arabidopsis* (*At*), tobacco (*Nt*), cucumber (*Cs*) and wild oat (*Af*). They are divided into three groups based on the number and type of the WRKY domains they contain. WRKY domains are black, putative basic nuclear localization signals are blue and leucine zippers are pink. Serine–threonine-rich regions are yellow, glutamine-rich regions are purple, proline-rich regions are green and acidic regions are red.

dimensional shapes of such complexes. Indeed, WRKY proteins might be part of multimeric protein–DNA complexes. Both WRKY protein-containing nuclear extracts and purified recombinant WRKYs from tobacco lose their DNA-binding activity when treated with the protein-dissociating agent deoxycholate<sup>18</sup>. Furthermore, some WRKY proteins contain potential leucine zippers (LZs), structures known to allow protein dimerization. They appear to be functional in *PcWRKY4* and *5*, because their deletion greatly reduces reporter-gene expression mediated by these proteins in yeast (R.S. Cormack *et al.*, unpublished).

#### Transcriptional regulation

ZAP1 and *PcWRKY1*, 4 and 5 can activate transcription in yeast<sup>10,12</sup> (R.S. Cormack *et al.*, unpublished), a feature that has been confirmed for ZAP1 and *PcWRKY1* in plant cells<sup>10,12</sup>. Although it has yet to be studied in detail, the primary structures of WRKY proteins have an abundance of potential transcriptional activation or repression domains (Fig. 2). A common feature of many domains affecting transcription is the predominance of certain amino acids,

Table 1. Identified members of the *Arabidopsis* WRKY superfamily of transcription factors

AtWRKY	Group	Chr.	EST <sup>a</sup>	Gene <sup>a</sup>	BAC.ORF	AtWRKY	Group	Chr.	EST <sup>a</sup>	Gene <sup>a</sup>	BAC.ORF
1	I	2	X92976/ZAP1 AI995838	AC007211 AC006955	F1013.1 F2818	19		4	AL049638 AL091613	F16J13.90 T8E18-end	
2	I	5	T44598 AA395490 N37131	AB026656	MXK23	20 21	I II	4 2	AA586133 T20410 AI992739	U93215 AQ010529	T15N24.90 T06B20.6 F24C8-end
3	I	2	T45479 AI099874 AI993164	AC006284 B77849 AL080571 AL096246	T4M8.23 T29F22-end F1G15-end T16K23-end	22 23	II II	4 2	T04811 F14417 F14438	AF007269 AC002337	IG002N01.6 T08I13.10
4	I	1	T22085 W43265 AA585810	AC007576	F7A19.5	24 25	I I	5 2	T42934	AB005233 AC002338 AC004165	MBK23 T9D9.6 T27E13.1
5		5	H77044 H77050 H77051 AI995170 AA605512	AB011485	MXH1.P3	26 27	I	5 5	AA585811 T22092 AI995443	AL093076 B09174 AB010697	T10P21-end T30A11-end M0J9.24
6	II	1	U75592 AA650675 H77127 AA394951 AA650826 AI992388	AC000375	F19K23.22	28 29 30 31 32 33	II II III II I I	4 4 5 4 4 2		AB009055 AL021713 AL035394 AB010696 AL022140 AL022198 AC004683	MXC20.3 T9A21.10 F9D19.20 MLE8.2 F1N20.170 F6I18.160 T19C21.4
7	II	4	N37775 T20578 R30038 AI992658	AC005861 AL078637	F23B24 T22A6.70	34 35 36	I II II	4 2 1		AC005499 AL022223 AC004238 AC010675	T6A23.33 M3E9.130 F19I3.6 T17F3.16
8	II	5		AB010698	MPL12.9					AC010852	T12P18
9			AI998645	AQ011596 B98122	F24A12TRC-end F24A12TRB-end	37 38	?	5	B23309	F28C5-end	
10	II	1		AC002328	F20N2	39	II	3	AB012244	MQJ16.9	
11	II	4	R64846 T88086 R30283 AI998936 T22071 T42669 Z29806 Z29805	AL080283	F3L17.120	40 41 42 43 44 45	II III II I	1 4 4 2 2	AC011713 AF080120 AL049876 AF076243 AC005397 AC005896	F23A5 F2P3.16 T22B4.50 T26N6.6 T3F17.22 F3G5.5	
								1/3?	ATU63815 AC010797	AT.I.24-4 IGF-F28J7	
12		2		AC003672	F16B22				AC011664 AC011624	F5A18 T18B3	
13			F14100						AC006526	F11C10.9	
14	II	4		AL078620	F23K16.40	46	III	2	AC006526	F11C10.9	
14		1	D88748 T20672	AC007060	T5I8.10	47 48	II	4 5	AF104919 AB023033	T15B16.12 K6M13	
15	II	2	Z25667 T04430 T43675 T21472 H36048 AI993841	AC002391	T20D16.5	49 50 51 52 53	II II III? III	4 5 5 5 4	AB017070 AC005965 AB019236 AB020744 AL078468	MNL12 T19G15 MXK3 K9E15 T32A16	
16	II	5	AA042185 AA395309	AB010693 B27842	K21C13.P1 T19B7-end	54 55	III III	2 2	AL035394 AC007660	F9D16.280 T7D17.7	
17		2	AA712348 R90490 AA067545 AI100579	AC006954	F25P17.13	56 57 58 59		1 3 2	AC007764 AL080748 AC008261 AC007019	F22C12 F1L14-end T4P13 F7D8.22	
18	II	4	U74179	AL031004 AL049607	F28M20.10 F11C18	60 61	II II	2 1	AC006585 AC011809	F27C12.8 F6A14	

<sup>a</sup>GenBank Accession no.

Abbreviations: Chr., chromosome; ORF, open reading frame; question marks denote either inconclusive group assignment or inconclusive chromosomal position.

including alanine, glutamine, proline, serine, threonine and charged amino acids<sup>21,22</sup>. At least two of the seven potential 'trans-regulatory' domains in *PcWRKY1* activate transcription in yeast<sup>12</sup>. The possibility that WRKY proteins possess both activator and repressor functions, as shown for the maize VP1 (Ref. 23), remains to be tested.

### Complexity of the WRKY family in *Arabidopsis*

The large amount of genomic and cDNA sequences available from *Arabidopsis* yields insights into the complexity of the WRKY family in a single plant species. In total, 61 distinct ORFs potentially encoding WRKY proteins can be found in the databases to date (Table 1). With the exception of *AtWRKY1*, which is identical to ZAP1 (Ref. 10), and *AtWRKY44*, which is defined by the *ttg2* mutant (C.S. Johnson and D.R. Smyth, pers. commun.), none of these proteins has been described before. We encourage the use of the designations used in Table 1 in future studies to avoid the confusion often caused when multiple names are assigned to a given gene member within large families.

The *AtWRKY* genes are randomly distributed over the five chromosomes and preliminary analyses suggest that they might all be present as single copies (I. Somssich and S. Robatzek, unpublished). Many of these putative WRKY proteins are represented by ESTs showing that the corresponding genes are expressed. By the number and sequence of their WRKY domains, these proteins can be assigned to the three major groups. Given that about two-thirds of the *Arabidopsis* genome has been sequenced to date, the total number of *WRKY* genes in this species might be as high as 100.

A phylogenetic tree of the *AtWRKY* proteins based on their WRKY domains (Fig. 3) clearly indicates that group II splits up into five distinct subgroups (IIa–e). The resulting refined classification is further substantiated by the presence of ten additional structural motifs that are conserved among subsets of *AtWRKY* family members. Each of these motifs occurs only in certain subgroups and each subgroup seems to be best defined by combinations of such motifs. In some cases, the sequences of these motifs can reveal clues about their potential functions. In addition to peptide sequences that might serve as nuclear localization signals<sup>24</sup>, a heptad repeat of bulky hydrophobic residues characteristic for LZs (Ref. 25) is present in some of the proteins. The heptad repeat occurs exclusively in members of subgroups IIa and IIb. Recent experiments have shown that the LZ region of *AtWRKY6* mediates dimerization (S. Robatzek and I.E. Somssich, unpublished).

An additional common feature that is found in the *WRKY* genes is the existence of an intron within the region encoding the C-terminal WRKY domain of group I members or the single WRKY domain of group II and III members. This intron position is highly conserved, being localized after the codon encoding arginine that is N terminal to the zinc-finger-like motif (Fig. 1). Strikingly, in all the genes encoding subgroup IIa and IIb members, the position of this intron is exactly 16 codons further towards the C terminus. In spite of the phylogenetic distance of their WRKY domains, members of all three groups have been shown to recognize W box elements, indicating that this is a general feature of the entire superfamily.

A few *AtWRKY* proteins do not fit neatly into any one (sub)group. For example, *AtWRKY10*, which carries only one WRKY domain, appears to be more related to group I (Fig. 3). This might be explained by the secondary loss of the N-terminal WRKY domain. Furthermore, based on the pattern of cysteine and histidine residues within their WRKY domains (Fig. 1), *AtWRKY38* and *AtWRKY52* could either belong to group III or represent members of a novel group (Fig. 3).

### Biological roles of WRKY factors

One of the most challenging questions concerns the regulatory processes governed by WRKY proteins. Clues might come partly from gene expression studies. Because many *WRKY* genes are themselves transcriptionally regulated, their distinct expression patterns might yield hints as to the regulatory functions of the encoded factors under particular biological conditions. In addition, a full understanding of the biological roles of these factors will require the identification of the target genes whose expression they affect.

#### Expression behavior of *WRKY* genes

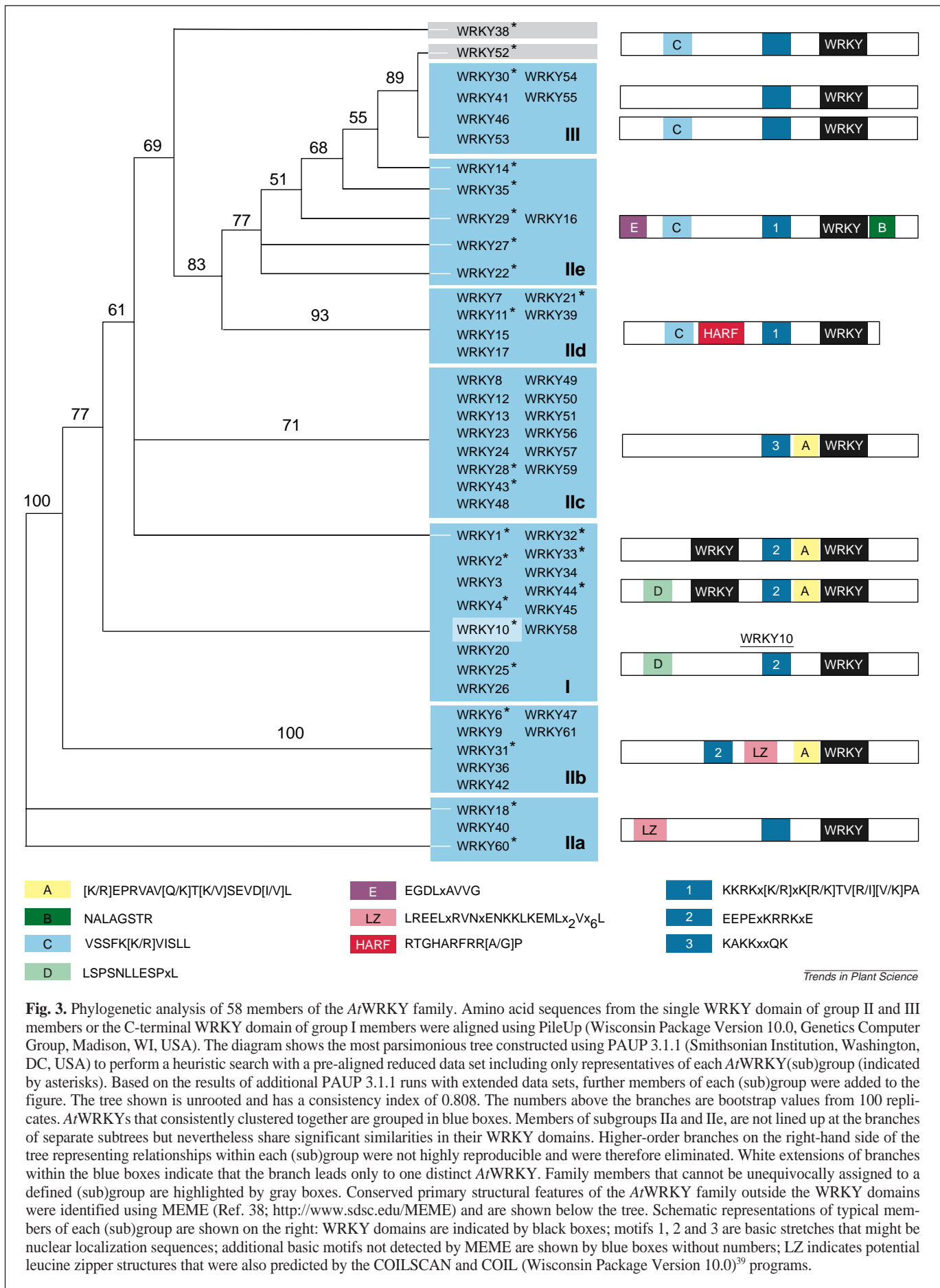
Current data point to many WRKY proteins having a regulatory function in the response to pathogen infection and other stresses. Effective plant defense against pathogenic microorganisms is associated with the concerted activation of a large variety of genes, occurring in several temporally distinct waves<sup>26</sup>. Increased levels of WRKY mRNA, protein and DNA-binding activity have been reported to be induced by infection with viruses<sup>19</sup>, bacteria (A. Dellagi and P. Birch, pers. commun.) or oomycetes<sup>12</sup>, by fungal elicitors<sup>9,20</sup> (R.S. Cormack *et al.*, unpublished), and by signaling substances such as salicylic acid<sup>18</sup>. In addition, *WRKY* gene expression has been shown to be upregulated in response to wounding<sup>11</sup> (S. Robatzek and I.E. Somssich, unpublished) and upon local mechanical stimulation of plant protoplasts<sup>27</sup>. Induced WRKY mRNA accumulation is often extremely rapid and transient, and seems not to require *de novo* synthesis of regulatory factors<sup>9,11</sup> (R.S. Cormack *et al.*, unpublished). This immediate-early expression behaviour indicates a role for the WRKY proteins in regulating subsequently activated secondary-response genes, whose products carry out the protective and defensive reactions.

Comparative expression studies with several *AtWRKY* genes also suggest that certain family members have a role in the regulation of senescence (S. Robatzek and I.E. Somssich, unpublished). Transcript levels of *AtWRKY4*, *6*, *7* and *11* are enhanced in senescent leaves. In transgenic *Arabidopsis* plants, an *AtWRKY6* promoter-*GUS* reporter gene is strongly activated in senescent leaves as well as in response to infection by pathogenic bacteria. As several genes are known to be highly expressed during both leaf senescence and defense, we might expect the existence of common regulatory mechanisms between these two physiological processes<sup>28</sup>.

Inspection of plant databases has revealed the existence of more than 500 WRKY ESTs identified from various tissue sources, including roots, leaves, inflorescences, abscission zones, seeds and vascular tissue, as well as from drought- or salt-stressed, or pathogen-infected tissue. Thus, *WRKY* genes appear to be expressed in numerous cell types and under different physiological conditions and could therefore participate in the control of a wide variety of biological processes.

#### Targets of *WRKY* regulation

As suggested by the general binding preference of WRKY proteins for W boxes, genes containing these promoter elements are likely targets of WRKY factors, and these include the *WRKY* genes themselves<sup>12</sup> as well as a large variety of defense-related genes of the *PR* type<sup>16,18</sup>. Additionally, gibberellic acid-induced expression of the wild-oat *α-Amy2/54* gene<sup>8</sup> and activation of the barley *HvLox1* gene in response to the defense and wound signaling molecule jasmonic acid<sup>29</sup> also appear to involve WRKY–W box interactions. Furthermore, a role has been suggested for SPF1 in the sucrose- or polygalacturonic-acid-induced expression of genes coding for sporamin and β-amylase in sweet potato<sup>7</sup>. However, as mentioned earlier, uncertainties about the



Trends in Plant Science

**Fig. 3.** Phylogenetic analysis of 58 members of the *AtWRKY* family. Amino acid sequences from the single WRKY domain of group II and III members or the C-terminal WRKY domain of group I members were aligned using PileUp (Wisconsin Package Version 10.0, Genetics Computer Group, Madison, WI, USA). The diagram shows the most parsimonious tree constructed using PAUP 3.1.1 (Smithsonian Institution, Washington, DC, USA) to perform a heuristic search with a pre-aligned reduced data set including only representatives of each *AtWRKY*(sub)group (indicated by asterisks). Based on the results of additional PAUP 3.1.1 runs with extended data sets, further members of each (sub)group were added to the figure. The tree shown is unrooted and has a consistency index of 0.808. The numbers above the branches are bootstrap values from 100 replicates. *AtWRKY*s that consistently clustered together are grouped in blue boxes. Members of subgroups IIa and IIe, are not lined up at the branches of separate subtrees but nevertheless share significant similarities in their WRKY domains. Higher-order branches on the right-hand side of the tree representing relationships within each (sub)group were not highly reproducible and were therefore eliminated. White extensions of branches within the blue boxes indicate that the branch leads only to one distinct *AtWRKY*. Family members that cannot be unequivocally assigned to a defined (sub)group are highlighted by gray boxes. Conserved primary structural features of the *AtWRKY* family outside the WRKY domains were identified using MEME (Ref. 38; <http://www.sdsc.edu/MEME>) and are shown below the tree. Schematic representations of typical members of each (sub)group are shown on the right: WRKY domains are indicated by black boxes; motifs 1, 2 and 3 are basic stretches that might be nuclear localization sequences; additional basic motifs not detected by MEME are shown by blue boxes without numbers; LZ indicates potential leucine zipper structures that were also predicted by the COILSCAN and COIL (Wisconsin Package Version 10.0)<sup>39</sup> programs.

exact binding site of SPFI means that more work is required to establish its role *in vivo*. To date, W boxes have been described as positive *cis*-acting elements upregulating transcription. However, in the case of the *Arabidopsis* *PRI* gene, the basal and salicylic acid-induced expression levels might be negatively regulated by W boxes<sup>17</sup>. SNI1, a negative regulator of *PR* gene expression, was recently identified in a genetic screening for second-site suppressors of the *Arabidopsis* mutation *npr1* (Refs 30,31). Interestingly, SNI1, which is nuclear localized, contains no obvious DNA-binding domain. One possible mode of SNI1 action would involve interaction with WRKY factors bound to the W box<sup>31</sup>.

The involvement of WRKY factors in regulating part of the defense program is further substantiated by a large-scale expression profiling study (J. Dangl and R.A. Dietrich, pers. commun.). Using a DNA microarray with 10 000 *Arabidopsis* ESTs, a group of 25 genes, including *PRI*, was identified whose expression responded coordinately to various pathogens as well as to other defense-inducing conditions. Within the first kilobase of their promoters, these genes shared only the W box motifs (TTGAC), with on average four copies, which were often clustered. By contrast, the promoters of a control set of genes not coordinately regulated with *PRI* contained, on average, less than two W boxes.

The only WRKY mutant so far described is *transparent testa glabra 2* (*ttg2*), which is based on a transposon insertion within *AtWRKY44/TTG2* (C.S. Johnson and D.R. Smyth, pers. commun.). In *ttg2*, the number of trichomes and their branching is reduced, as is anthocyanin pigmentation of the seed coat, together with a loss of mucilage. This pleiotropic phenotype resembles that of *ttg1*, which is defective for a regulatory protein of the WD40-repeat type<sup>32</sup>. *AtWRKY44/TTG2* and *TTG1* might therefore act in the same regulatory cascade, controlling a common set of genes. The extensive use of reverse genetics to obtain additional tagged WRKY mutants, as well as the generation of WRKY promoter-reporter gene and WRKY overexpressor lines, will allow us to gain a more comprehensive understanding of the various biological roles of WRKY proteins. Furthermore, inducible expression systems<sup>33</sup> could be used for controlled temporal overexpression of WRKY transgenes in their own loss-of-function background; combined with methods of large-scale gene expression profiling (e.g. differential display, DNA chips), this should facilitate the identification of defined WRKY target genes. In a similar way, the *Arabidopsis* *NAP* gene was identified as a target of the APETALA3-PISTILLATA transcription factor dimer<sup>34</sup>.

## Conclusions

WRKY proteins have only recently been identified as a new family of transcription factors. In *Arabidopsis*, this family appears to be nearly as complex as the well-known MYB family<sup>35</sup>, but it is restricted to the plant kingdom. This suggests that WRKY genes originated concurrently with the major plant phyla. Current information suggests that WRKY factors play a key role in regulating the pathogen-induced defense program. The exposure of plants to a wide variety of biotic or abiotic stresses connected with their sessile, autotrophic lifestyle could be one major factor in the enormous expansion of the WRKY family during evolution. In addition, the extensive metabolic changes associated with the establishment of defense responses<sup>26</sup> or senescence<sup>36</sup> might require an elaborate regulatory system.

WRKY proteins also seem to be involved in other plant-specific processes, such as trichome development and the biosynthesis of secondary metabolites. Thus, they appear to participate in controlling the expression of a plethora of genes. As with other large gene families, the problem of functional redundancy will complicate genetic attempts to determine the role of individual

WRKY proteins. Comparative studies in lower plants (e.g. ferns, mosses and algae) can give clues to whether WRKY gene diversification correlates with increasing developmental and metabolic pathway complexity. Furthermore, generating *Arabidopsis* knock-out lines that affect several members of individual subgroups might help to 'wrky' matters out.

## Acknowledgements

We thank Hiroshi Sano (NAIST, Japan); David R. Smyth (Monash University, Australia); Zhixiang Chen (University of Idaho, USA); Jeff Dangl (University of North Carolina, USA); Robert Dietrich (Novartis, Research Triangle, USA); Alia Dellagi and Paul Birch (Scottish Crop Research Institute, UK), for providing preprints of unpublished data; and Klaus Hahlbrock for critical reading of the manuscript and continuous support.

## References

- 1 Bevan, M. *et al.* (1998) Analysis of 1.9 Mb of contiguous sequence from chromosome 4 of *Arabidopsis thaliana*. *Nature* 391, 485–488
- 2 Mewes, H.W. *et al.* (1997) Overview of the yeast genome. *Nature* 387, 7–8
- 3 Clarke, N.D. and Berg, J.M. (1998) Zinc fingers in *Caenorhabditis elegans*: finding families and probing pathways. *Science* 282, 2018–2022
- 4 Laudet, V. *et al.* (1992) Evolution of the nuclear receptor gene superfamily. *EMBO J.* 11, 1003–1013
- 5 Gellon, G. and McGinnis, W. (1998) Shaping animal body plans in developmental and evolution by modulation of *Hox* expression patterns. *BioEssays* 20, 116–125
- 6 Riechmann, J.L. and Meyerowitz, E.M. (1997) MADS domain proteins in plant development. *Biol. Chem.* 378, 1079–1101
- 7 Ishiguro, S. and Nakamura, K. (1994) Characterization of a cDNA encoding a novel DNA-binding protein, SPFI, that recognizes SP8 sequences in the 5' upstream regions of genes coding for sporamin and  $\beta$ -amylase from sweet potato. *Mol. Gen. Genet.* 244, 563–571
- 8 Rushton, P.J. *et al.* (1995) Members of a new family of DNA-binding proteins bind to a conserved *cis*-element in the promoters of  $\alpha$ -Amy2 genes. *Plant Mol. Biol.* 29, 691–702
- 9 Rushton, P.J. *et al.* (1996) Interaction of elicitor-induced DNA binding proteins with elicitor response elements in the promoters of parsley PR1 genes. *EMBO J.* 15, 5690–5700
- 10 de Pater, S. *et al.* (1996) Characterization of a zinc-dependent transcriptional activator from *Arabidopsis*. *Nucleic Acids Res.* 24, 4624–4631
- 11 Hara, K. *et al.* (2000) Rapid systemic accumulation of transcripts encoding a tobacco WRKY transcription factor upon wounding. *Mol. Gen. Genet.* 263, 30–37
- 12 Eulgem, T. *et al.* (1999) Early nuclear events in plant defense: rapid gene activation by WRKY transcription factors. *EMBO J.* 18, 4689–4699
- 13 Berg, J.M. and Shi, Y. (1996) The galvanization of biology: a growing appreciation for the roles of zinc. *Science* 271, 1081–1085
- 14 Mackay, J.P. and Crossley, M. (1998) Zinc fingers are sticking together. *Trends Biochem. Sci.* 23, 1–4
- 15 Fukuda, Y. and Shinshi, H. (1994) Characterization of a novel *cis*-acting element that is responsive to a fungal elicitor in the promoter of a tobacco class I chitinase gene. *Plant Mol. Biol.* 24, 485–493
- 16 Rushton, P.J. and Somssich, I.E. (1998) Transcriptional control of plant genes responsive to pathogens. *Curr. Opin. Plant Biol.* 1, 311–315
- 17 Lebel, E. *et al.* (1998) Functional analysis of the regulatory sequences controlling *PR-1* gene expression in *Arabidopsis*. *Plant J.* 16, 223–233
- 18 Yang, P. *et al.* (1999) A pathogen- and salicylic acid-induced WRKY DNA-binding activity recognizes the elicitor response element of tobacco class I chitinase gene promoter. *Plant J.* 18, 141–149
- 19 Wang, Z. *et al.* (1998) An oligo selection procedure for identification of sequence-specific DNA-binding activities associated with plant defense. *Plant J.* 16, 515–522

- 20 Fukuda, Y. (1997) Interaction of tobacco nuclear proteins with an elicitor-responsive element in the promoter of a basic class I chitinase gene. *Plant Mol. Biol.* 34, 81–87
- 21 Triezenberg, S.J. (1995) Structure and function of transcriptional activation domains. *Curr. Opin. Genet. Dev.* 5, 190–196
- 22 Hanna-Rose, W. and Hansen, U. (1996) Active repression mechanisms of eukaryotic transcription repressors. *Trends Genet.* 12, 229–234
- 23 Hoecker, U. *et al.* (1995) Integrated control of seed maturation and germination programs by activator and repressor functions of Viviparous-1 of maize. *Genes Dev.* 9, 2459–2469
- 24 Garcia-Bustos, J. *et al.* (1991) Nuclear protein localization. *Biochim. Biophys. Acta* 1071, 83–101
- 25 Landschulz, W.H. *et al.* (1988) The leucine zipper: a hypothetical structure common to a new class of DNA-binding proteins. *Science* 240, 1759–1764
- 26 Somssich, I.E. and Hahlbrock, K. (1998) Pathogen defense in plants – a paradigm of biological complexity. *Trends Plant Sci.* 3, 86–90
- 27 Gus-Mayer, S. *et al.* (1998) Local mechanical stimulation induces components of the pathogen defense response in parsley. *Proc. Natl. Acad. Sci. U. S. A.* 95, 8398–8403
- 28 Quirino, B.F. *et al.* (1999) Diverse range of gene activity during *Arabidopsis thaliana* leaf senescence includes pathogen-independent induction of defense-related genes. *Plant Mol. Biol.* 40, 267–278
- 29 Rouster, J. *et al.* (1997) Identification of a methyl-jasmonate-responsive region in the promoter of a lipoxygenase-1 gene expressed in barley grain. *Plant J.* 11, 513–523
- 30 Cao, H. *et al.* (1997) The *Arabidopsis NPR1* gene that controls systemic acquired resistance encodes a novel protein containing ankyrin repeats. *Cell* 88, 57–63
- 31 Li, X. *et al.* (1999) Identification and cloning of a negative regulator of systemic acquired resistance, SN11, through a screen for suppressors of *npr1-1*. *Cell* 98, 329–339
- 32 Walker, A.R. *et al.* (1999) The *TRANSPARENT TESTA GLABRA1* locus, which regulates trichome differentiation and anthocyanin biosynthesis in *Arabidopsis*, encodes a WD40 repeat protein. *Plant Cell* 11, 1337–1349
- 33 Gatz, C. and Lenk, I. (1998) Promoters that respond to chemical inducers. *Trends Plant Sci.* 3, 352–358
- 34 Sablowski, R.W.M. and Meyerowitz, E.M. (1998) A homolog of *NO APICAL MERISTEM* is an immediate target of the floral homeotic genes *APETALA3/PISTILLATA*. *Cell* 92, 93–103
- 35 Martin, C. and Paz-Ares, J. (1997) MYB transcription factors in plants. *Trends Genet.* 13, 67–73
- 36 Gan, S. and Amasino, R.M. (1997) Making sense of senescence. *Plant Physiol.* 113, 313–319
- 37 Altschul, S.F. *et al.* (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res.* 25, 3389–3402
- 38 Bailey, T.L. and Elkan, C. (1994) Fitting a mixture model by expectation maximization to discover motifs in biopolymers. In *Proceedings of the Second International Conference on Intelligent Systems for Molecular Biology* (Altmann, R., ed.), pp. 28–36, AAAI Press
- 39 Lupas, A. (1996) Coiled coils: new structures and new functions. *Trends Biochem. Sci.* 21, 375–382

Thomas Eulgem, Paul Rushton, Silke Robatzek and Imre Somssich\* are at the Max-Planck-Institut für Züchtungsforschung, Abteilung Biochemie, Carl-von-Linné-Weg 10, D-50829 Köln, Germany; Thomas Eulgem is currently at the Dept of Biology, 108 Coker Hall CB#3280, University of North Carolina, Chapel Hill, NC 27599-3280, USA.

\*Author for correspondence (tel +49 221 5062310; fax +49 221 5062313; e-mail somssich@mpiz-koeln.mpg.de).

# Plant one-carbon metabolism and its engineering

Andrew D. Hanson, Douglas A. Gage and Yair Shachar-Hill

**The metabolism of one-carbon (C<sub>1</sub>) units is vital to plants. It involves unique enzymes and takes place in four subcellular compartments. Plant C<sub>1</sub> biochemistry has remained relatively unexplored, partly because of the low abundance or the lability of many of its enzymes and intermediates. Fortunately, DNA sequence databases now make it easier to characterize known C<sub>1</sub> enzymes and to discover new ones, to identify pathways that might carry high C<sub>1</sub> fluxes, and to use engineering to redirect C<sub>1</sub> fluxes and to understand their control better.**

One-carbon (C<sub>1</sub>) metabolism is essential to all organisms. In plants, it supplies the C<sub>1</sub> units needed to synthesize proteins, nucleic acids, pantothenate and many methylated molecules<sup>1</sup>. Fluxes through C<sub>1</sub> pathways are particularly high in plants that are rich in methylated compounds such as lignin, alkaloids and betaines because methyl moieties make up several percent of their dry weight<sup>2</sup>. Transfers of C<sub>1</sub> units are also central to the massive photorespiratory fluxes that occur in all C<sub>3</sub> plants<sup>3</sup>. In spite of the fundamental significance of these roles, and the interest in the metabolic engineering of lignin<sup>2</sup>, betaines<sup>4</sup> and photorespiration<sup>3</sup>, there is much that is not understood about the enzymes, pathways and regulatory mechanisms of plant C<sub>1</sub> metabolism. In part this is because of the obstacles that C<sub>1</sub> metabolism presents for classical biochemistry and genetics: its enzymes can be of low abundance and/or exist as

several isoforms, mutants are lacking, and its key intermediates – C<sub>1</sub> substituted folates – are labile and hard to quantify. Fortunately, classical approaches to C<sub>1</sub> metabolism can now be complemented by genomics-driven approaches that exploit the fast-growing DNA sequence databases. Accordingly, this review has three aims:

- To illustrate how genomics-based approaches are advancing our knowledge of plant C<sub>1</sub> biochemistry.
  - To bring together biochemical and genomics-derived data to show which C<sub>1</sub> pathways might operate in plants, and where they operate in the cell.
  - To examine progress towards engineering C<sub>1</sub> metabolism.
- Nucleotide sequence information – from genomes, cDNAs and ESTs – can be used to complement biochemical approaches in several ways. Because most enzymes of C<sub>1</sub> metabolism are highly