

Microbial Relatives of Seed Storage Proteins: Conservation of Motifs in a Functionally Diverse Superfamily of Enzymes

Jim M. Dunwell,¹ Paul J. Gane^{2,*}

¹ Department of Agricultural Botany, School of Plant Sciences, The University of Reading, Whiteknights, PO Box 221, Reading RG6 6AS, United Kingdom

² Institute of Food Research, Reading Laboratory, Earley Gate, Whiteknights Road, Reading RG6 6BZ, United Kingdom

Received: 25 April 1997 / Accepted: 29 July 1997

Abstract. Plant storage proteins comprise a major part of the human diet. Sequence analysis has revealed that these proteins probably share a common ancestor with a fungal oxalate decarboxylase and/or related bacterial genes. Additionally, all these proteins share a central core sequence with several other functionally diverse enzymes and binding proteins, many of which are associated with synthesis of the extracellular matrix during sporulation/encystment. A possible prokaryotic relative of this sequence is a bacterial protein (SASP) known to bind to DNA and thereby protect spores from extreme environmental conditions. This ability to maintain cell viability during periods of dehydration in spores and seeds may relate to absolute conservation of residues involved in structure determination.

Key words: Seed storage proteins — Enzyme superfamily — Protein domain — Germin — Oxalate oxidase — Histidine cluster — Mannose metabolism

Introduction

Plant storage proteins, particularly those found in seeds, form a major part of the human diet. Many such proteins

(including the vicilins, legumins, and globulins) are now known to have two domains, each similar in sequence to a family of proteins (Bäumlein et al. 1995) the best known of which is germin, the predominant protein produced during the early phase of germination of the wheat embryo (Lane et al. 1992). Germin is a glycosylated, homopentameric protein with exceptional resistance to proteases and hydrogen peroxide, the latter feature related to its function as an oxalate oxidase (Lane et al. 1993) that generates peroxide. There are several germin-like proteins found in dicotyledonous species (Heintzen et al. 1994; Ono et al. 1996) (see also the multigene family in *Arabidopsis*, accessions U75187–U75207, U95034–U95036) and gymnosperms (Domon et al. 1995). The function of these proteins is unknown; most do not have any oxalate oxidase activity, although at an amino acid level there is a high level of similarity to the cereal ox-ox enzymes (for a recent discussion see Berna and Bernier 1997).

The present study started from the specific identification (Lane et al. 1991) of the so-called “germin box” (HI/THPRATEI), a conserved sequence shared by the germins and spherulins—the latter are a group of proteins produced in the slime mold *Physarum polycephalum* during encystment (Bernier et al. 1987). Previous PROSITE analysis had identified at PDOC00597 a germin family signature PS00725 which comprised the germin box. Additionally, there is a three-element PRINTS fingerprint GERMIN which is based on alignment of 12 sequences and a ProDom domain 2426 (ProDom release 34.1) found in 10 proteins. In each case the analysis is

* Present address: The Drug Design Group, Department of Pharmacology, University of Cambridge, Tennis Court Road, Cambridge CB2 1QJ, United Kingdom

Correspondence to: J.M. Dunwell; e-mail J.Dunwell@reading.ac.uk

severely outdated since it is based on only a small subset of those sequences now available. (The SwissProt database contains only a minority of the sequences available.)

To refine these analyses, various methods were adopted in an attempt to identify possible progenitors of the germin/germin-like proteins, and also to assess the occurrence of any conserved motifs with potential active site function (see Gane et al. 1997).

Methods

Initially, BLAST searches were conducted on various germin sequences in the GenBank database, and these were linked to BLOCKS (Henikoff and Henikoff 1991), MEME (Bailey and Elkan 1994), ProDom, and PROSITE analyses. Taken together, these analyses showed (1) that only certain residues within the germin box are absolutely conserved across all proteins and (2) that other regions of the gene are equally conserved. Specifically, the germin box is part of a 20/21-AA motif G(X)₅HXH(X)₁₁G (part of PRINTS motif GERMIN1) which is followed (usually after 15 residues) by a second motif of 16 AAs, G(X)₅P(X)₄H(X)₃N (part of PRINTS motif GERMIN2). When one uses the bean storage protein phaseolin (Lawrence et al. 1994) and the related sucrose-binding protein (Braun et al. 1996) as structural references, these two histidine-containing motifs are part of the β -strands designated, respectively, C/D and G/H within the two β -barrel elements. Considering the first motif, the two flanking glycines correspond to the strictly conserved Gly²⁴⁹ and Gly²⁶⁹ which assist in the formation of the short interstrand loops B-C and D-E within the C-terminal β -barrel of phaseolin. Similarly, the second motif contains a proline corresponding to conserved Pro³⁰³ which is part of the interstrand loop between strands G and H. The variable space between motifs is equivalent to the insertions (seven to 25 residues) tolerated in the E-F loop (see Gane et al. 1997).

Gene families found to contain these two specific motifs are considered below, with summary details provided in Table 1 and Fig. 1. In order to provide a link to previous publications (Bäumlein et al. 1995; Braun et al. 1996; Lane et al. 1991; Lawrence et al. 1994) these summaries include phaseolin (gi|230247), a spherulin (P09351), a vicilin (Z54364) found in spores of the Ostrich fern *Matteuccia struthiopteris* (the C-terminal domain of this sequence is particularly similar to wheat germin) and the reference standard wheat germin gf2.8 (P15290).

It should be noted that there are two examples described below (i.e., VIRT18179 and VIRT13653) in which TBLASTN searches revealed previously unknown coding regions in upstream sequences of other genes. The explanation for the presence of these cryptic genes (sometimes encoded by sequences in different frames) probably lies in cloning artifacts caused by accidental ligation of two different DNA fragments.

Results

Oxalate Decarboxylase and Other Duplicated Proteins

This study revealed three examples of two-domain proteins (Fig. 2) in which each domain contains both of the two motifs, and there is significant similarity to GLPs and their relatives. The first, and best described in terms of function, is the oxalate decarboxylase (ODC) enzyme (I25120) (Mehta and Datta 1991) from the wood-rotting

fungus *Collybia velutipes* (Fig. 2A). It is of especial interest that this oxalate-degrading enzyme should be a duplicated version of the only other enzyme with a related function—namely, the cereal oxalate oxidase described above. (However, it does not seem to be related to the oxalyl-CoA decarboxylase enzyme, gi|1086099, from *Oxalobacter formigenes*.) The closest neighbor [WU-BLASTP, P(N)3.0e-55; 37% identity and 54% similarity over 326 AA] to the *Collybia* ODC sequence is the second example, a hypothetical protein (accession D90907) from the cyanobacterium *Synechocystis* (Kaneko et al. 1996) (Fig. 2B) which is shown in a WU-BLASTP search as being most similar [P(N)2.5e-216] to the *Arabidopsis* GLP10 (U95036); it is 30% identical and 49% similar over a distance of 139 AA. The third example of duplication (accession P42106) is another protein (Yoshida et al. 1995) of unknown function, from *Bacillus subtilis*, and was identified in the ProDom study referred to later. It is one of 17 proteins containing the 52-AA domain 1428 (ProDom release 34.1) which spans the two motifs identified here, and it has as its closest neighbor [P(N)0.00028] the polyketide synthase P23157 referred to below.

Despite the fact that the critical residues within each motif have been conserved in all three examples, there are a number of significant differences between the sequences of these proteins. First, the gap between motifs in P42106 is 15 AA (the overall minimum), and 20 AA in the other two examples. Second, overlaps of the duplicated regions are of different length, and in addition, BLASTP searches with BEAUTY annotation (Worley et al. 1990) show that the stretches of sequence showing significant similarity to other proteins are located primarily in the C-terminal section of oxalate decarboxylase, whereas they are located mostly in the N-terminal half of accession P42106. The third protein, D90907, has equal similarity (particularly to germins/vicilins) in each half.

In terms of their evolution, there are a number of possible origins for this duplication, which previously has never been found outside the plant kingdom. The simplest hypothesis is that the two domains of each protein have diverged differentially after a single duplication event. Alternatively, duplication might have occurred three times, most recently in *Synechocystis* (D90907), in which the two domains show equal similarity to other proteins. The third possibility is that this latter protein is the only one to be a product of duplication, and the other two examples have evolved from a homologous recombination event between two similar DNA sequences. In this case, D90907 would be the direct and sole ancestor of the higher plant storage proteins.

Auxin-Binding Proteins

An extremely high level of similarity (66% identity; 77% similarity) was identified between two auxin-binding

Table 1. Summary of single- and double-domain proteins which contain motifs 1 and 2^a

Species	Accession	Length (AA)	Name/function
i) Single domain			
Prokaryote:Archaea			
<i>Desulfurococcus</i> sp.	D84067	118	VIRT18179, unknown
<i>Methanococcus jannaschii</i>	U67602	125	?PMI/GDP
Prokaryote:Eubacteria:Proteobacteria			
<i>Erwinia chrysanthemi</i>	Q05527	110	Unknown
<i>Escherichia coli</i>	P38522	121	Aldehyde dehydrogenase
<i>Erwinia chrysanthemi</i>	L39897	132	ORF 1/unknown
<i>Enterobacter aerogenes</i>	U60777	140	Pep1
<i>Desulfibrio desulfuricans</i>	Z11975	163	VIRT13653, unknown
<i>Neisseria meningitidis</i>	L09188	180	Deoxyglucose epimerase
<i>Escherichia coli</i>	gi 1788802	233	Unknown
<i>Escherichia coli</i>	P17410	280	Cel operon repressor
<i>Yersinia enterocolitica</i>	U46859	465	Mannose pyrophosphorylase
<i>Xanthomonas campestris</i>	P29956	466	Phosphomannose isomerase
Prokaryote:Eubacteria:Firmicutes:Actinomycete			
<i>Streptomyces coelicolor</i>	U37580	77	Membrane-spanning protein
<i>Mycobacterium tuberculosis</i>	Z81360	116	Unknown
<i>Streptomyces cyaneus</i>	Q02586	154	CurC/?cyclase
<i>Mycobacterium tuberculosis</i>	gi 1781124	263	?regulatory protein
Prokaryote:Eubacteria:Firmicutes:LowG+C gram positive			
<i>Bacillus subtilis</i>	gi 1881251	113	YdbB/unknown
<i>Bacillus subtilis</i>	P54430	186	YrkC/unknown
<i>Acholeplasma laidlawii</i>	S33518	369	Unknown
<i>Staphylococcus aureus</i>	U81973	371	Cap5F/unknown
<i>Bacillus subtilis</i>	P39631	432	Spore polysacch. synth.
<i>Clostridium thermocellum</i>	P26208	448	β-glucosidase A
Prokaryote:Eubacteria:Cyanobacteria			
<i>Synechocystis</i> sp. PCC6803	D64001	128	Phosphomannose isomerase
<i>Synechocystis</i> sp. PCC6803	D90909	135	Unknown
<i>Synechocystis</i> sp. PCC6803	D90910	143	Unknown/?cytochrome c551
<i>Synechococcus</i> sp.	S04426	150	Unknown/?PMI
Eukaryote:Fungi:Ascomycetes			
<i>Saccharomyces cerevisiae</i>	P47096	177	?3-HAO
<i>Saccharomyces cerevisiae</i>	S53039	179	Unknown
Eukaryote:Myxomycetes			
<i>Physarum polycephalum</i>	P09351	248	Spherulin
Eukaryote:Plant:Angiosperm			
<i>Arabidopsis thaliana</i>	P33487	198	Auxin-binding protein
<i>Triticum aestivum</i>	P15290	201	Germin/oxalate oxidase
<i>Prunus persica</i>	gi 1916807	214	Auxin-binding protein
<i>Arabidopsis thaliana</i>	Q05212	230	DRT 102/DNA damage repair
<i>Vicia faba</i>	X95995	1641	ENBP1/zinc finger protein
Eukaryote:Animal			
<i>Caenorhabditis elegans</i>	gi 1082118	104	Unknown
<i>Caenorhabditis elegans</i>	gi 1707132	159	Unknown
<i>Caenorhabditis elegans</i>	P39645	190	RFBC/epimerase
<i>Caenorhabditis elegans</i>	Z70755	207	3-HAO
<i>Caenorhabditis elegans</i>	gi 726397	221	Unknown
<i>Caenorhabditis elegans</i>	gi 1572766	349	Unknown
<i>Caenorhabditis elegans</i>	Z50070	645	Unknown
<i>Caenorhabditis elegans</i>	gi 1572765	910	Unknown
<i>Caenorhabditis elegans</i>	U00043	1070	Unknown
<i>Rattus norvegicus</i>	P21816	200	Cysteine dioxygenase
<i>Rattus rattus</i>	D44494	286	3-HAO
<i>Rattus norvegicus</i>	X59993	1214	Zinc finger protein
ii) Double domain			
Prokaryote:Eubacteria:Firmicutes:LowG+C gram positive			
<i>Bacillus subtilis</i>	P42106	337	Unknown
Prokaryote:Eubacteria:Cyanobacteria			
<i>Synechocystis</i> sp. PCC6803	D90907	394	Unknown
Eukaryote:Fungi:Basidiomycete			
<i>Collybia velutipes</i>	I25120	447	Oxalate decarboxylase
Eukaryote:Plant:Pteridophyte			
<i>Matteuccia struthiopteris</i>	Z54364	504	Spore vicilin
Eukaryote:Plant:Angiosperm			
<i>Phaseolus vulgaris</i>	gi 230247	397	Phaseolin

Accession

Motif 1

AA

Motif 2

(i)

u37580	T T D S Q K P H A	Q D E V Y F V V S	15
z11975	E K I S A H T S T	G D A F V L A L E G	15
s04426	Q Q L S L Q R H Q	Q R Q H W L V V Q G	15
u46859	Q R T A T Q I H H	H R A H W V V V S	15
d64001	H R L S L Q M H H	H R S H W I V V S	15
p29956	A T L S L Q M H H	H R A H W I V V S	15
u81973	I T K G N W H H	T K N K F L V V S	20
s33518	I T K G N W H H	T K N K F L V V S	20
q05527	A I G T P H K H D	I H D Q I A Y A A	15
p38522	T T T G E R I K H	Q G E I G T L E	15
u67602	S K T L L K Y L	T S E I Y Y I L E	15
d84067	Q T V K K H Y L	H Q Y L F Y I M S	15
z81360	G T A E P A P T R	E E T V Y V L D	23
l39897	A Q A E P H Y V P	D Y E T A I Y L L K	18
gi1881251	E Y D W H H V	D S D L F I L E	16
x95995	N D C E S M H Y D N	V Q D R C S S Q	138
d90909	Q Q I V Q R H S H P	D G Q D T W V M L K	16
u60777	Q G I N R R L Y H P	A A V T F V L S	16
p54430	E D I G L E I H P	N V D Q F L R I E Q	21
d90910	A E T G W H S H P	V P S F G I L L E	16
q02586	E R I S E H Y H P	Y S E F V Y V E	15
p26208	D V A C D H Y H R	Y E E D I K I M K E I	
p33487	S E T P I R H S H P	C E E V F V L K	24
p09351	I N L P T H P	R A T I N F I A K	23
p15290	G T N P P H I H P	R A T I G I V M K	23
gi1781124	G A R I E R H H P	S H O I V Y P S A	15
q05212	S V E P A H H T	F G H D L V I K	17
p17410	E S I S G L Q H D	Y Y E F T L V L E	15
gi1788802	Q W E N A F F P W T	L N D E I D M V L E	16
p39631	V I K A F H Y E	K Q D D L W F F P T	24
s53039	A T F Y Q E L H E	D E E I R Y C L E	22
l33181	N V I R G M F Q M P A E H D K L V Y C V N		27
l09188	V L R G L Y Q T	E N T Q G K L V R V V	27
p47096	P N E R T G Y H I	N P T P W F Y Q K K	23
p40034	D A Y T D F H L D F	A G T S V Y Y N I S	45

G	S V V Y V P A G V A K F H H
G G	E S I I M P A G Q P H S V S A
G	O S L D A I G E W H R L Q A
N	E S T Y A V G V A H S I E N
N	Q S T Y V Q C T A H R L E N
N	O S T Y P L G V T H R L K N
G D K L	E V V D P V G Y T H N I E N
G G G	E K L E V V D P P G Y T H N I E N
G	D A Y M A V K N E M H G V V S
G	Q S Y A N T G I P H S F S N
G	D T I Y P K T P H K I E N
G	D I F L V K P K T V H W V I N
G	E L R S H N R T D R Q A L L
G	E F L Y P K G T V H Q P R N
N	D S L L P K G P T V H R T R S
G	E A V F P A G C P H Q V R N
G	E V A I A E K N Q V H G A I N
D	D S L L V P A G T T H S H W N
G	S A I V P A G T W H N V I N
G	D A I A E V V N T V H N G R N
D	Q G L M P I D M R H R F R N

(ii)

u39645	V L R G L H T Q P	H N G K L V T V V S	25
x59993	T T N L L D V	S D A A N V M V Y V	103
gi1707132	E Q F Y E P Q V Q K	E D V I S L V E	20
p34650	S I Y Q L P Y S E	S C S V L T L Y	17
p34949	S V T E Y K V L A	L D S A S I L L M V Q	17
p21816	H G S S I D H T	D S H C F L K L L Q	19
z70755	P N Q R K D F H L	E E G E F F F Q R K	19
d44494	P N T R K D Y H I	E E G E F V F Y Q L E	19
u00043	N C Y T D F I D F	S G T S V W Y H V L K	26
z50070	N S Y T D F V D F	G G T S V Y F H V F K	45
gi1572765	S Y T D F H V D F	G G S S V Y Y H I L K	45
gi1572766	R S G T A I H I D P	L G T S A W N S L L Q	59

G D N K	H A F W P A G F L H G F Q V
G G	D V V F P A G A P H Q V H N
G	D L I V P K G L S H R F T T
G	E V V F P G A T H D A E R
G	G V L F P G A N E S V S L

G	E M F M L P A R V E H S P Q R
G	E I F L L P A R V P H S P Q R
G	D T M L P S G W I A V Y T
G	Q T L L P A G W I A V L T
G	Q T L L P A G W I A V L T
G	E T M F V P S G W W H V V I N

(iii)

p42106a	D A F P L V H K	D T H G I L V L D	15
p42106b	D R I V D Y H E	Y H T T F Y C L E	15
d90907a	A I R E L W H A	N A A M W A Y V M E	20
d90907b	A M R Q L W G P	N A D W Q Y V L D	20
i25120a	A I R E L W H K	N A W A Y V L K	20
i25120b	A L R E L W H P	T E D W T F F I S	20
z54364b	A V L A P W N P	R A T I A L T K	21
230247b	A L F V P H Y S	K A I V I L V N E	27

G	D Y A N I P A G T P S Y R M
G	D F L H V P A N T V H S Y R L
G	G L W Y F P R G W G H S I E G
G	D V G Y V P K G Y G H A I R N
G	D L W Y F P G I P H S L Q A
G	D I A Y V P A S M G H Y V E N
G	S V F F V P Q N F P M C Q A
D	D V F V P A A Y P V A I K A

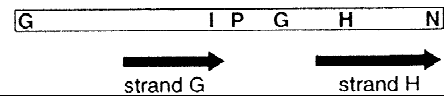
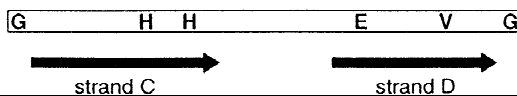


Fig. 1. Details of the two conserved motifs and their spacing, present in a range of proteins identified by their GenBank accession numbers. Section (i) denotes microbial and plant proteins, section (ii) animal proteins, and section (iii) two-domain proteins. Suffixes (a) and (b) represent, respectively, the first and second domains in these two-domain proteins. Highly conserved residues are shaded and are shown in the boxes below each alignment. Nomenclature of β -strands is according to Lawrence et al. (1994).

proteins (ABPs) from *Prunus persica* (gi|1916807, gi|1916809) and several germins (e.g., U75193) from *Arabidopsis thaliana*. Considering only the two motifs, this similarity to U75193 increases to 17/21 (81%) identity for motif 1 and 14/16 (88%) identity for motif 2.

Indeed, on the basis of their sequence, these peach ABPs should probably be reclassified as germin-like proteins (GLPs). However, extensive similarities to GLPs were also detected in sequences from functionally defined ABPs (Venis and Napier 1995) such as that from *A.*

and Bernier 1997). It is therefore possible that an ABP could modify the expression of germin by reducing the level of free auxin.

Polyketide Synthases

Amongst the sequences with the closest similarity [G(X)₅HYHPYSEE(X)₆G] to GLPs is accession M33704. This *Streptomyces cyaneus* gene (Bergh and Uhlen 1992), *curC* (probably a cyclase), is part of the synthetic pathway of the antibiotic curamycin. Closely related genes in other *Streptomyces* species include P23157, Q05362, and P16558 (see Fig. 4 in Blanco et al. 1993, which also identified residues 47–60, PGERISEHYHPYSE, cf. motif 1, as being conserved amongst several β-glucosidases; see P26208 from *Clostridium thermocellum*). Additional members include the *B. subtilis* sequence P54430 and the sequence Pep1 (Smith et al. 1993) (U60777) from *Enterobacter aerogenes*. The smallest (110–132 AA) proteins in this group are those from *Erwinia chrysanthemi* (accessions Q05527, L39897) and *B. subtilis* (gi|1881251).

Phosphomannose Isomerase and Related Enzymes Concerned With Polysaccharide Synthesis

Phosphomannose isomerases (PMIs) are zinc-containing enzymes that catalyze the interconversion of mannose-6-phosphate and fructose-6-phosphate. The class II PMIs (Proudfoot et al. 1994) are involved in a variety of pathways including capsular polysaccharide biosynthesis and D-mannose metabolism. Their similarity to proteins of the previous section was already known from the ProDom database, which includes 17 proteins with the domain 1428. These 17 include PMIs, polyketide synthases, and the duplicated *B. subtilis* sequence referred to above. Notable amongst the former category are the spore polysaccharide protein (P39631) from *B. subtilis*, the Cap5F and Cap8F proteins from *Staphylococcus aureus* (Sau and Lee 1996), and the closely related (>90%) sequence (S33518) from the mycoplasma *Acholeplasma laidlawii*. This class of PMI includes, toward their C-terminus, histidine-containing motifs similar in sequence and spacing to those described above. One difference is that the histidines of the first motif are often displaced to give a G(X)₇HXH(X)₆G sequence. Of this group of proteins, the smallest (128 AA) is from *Synechocystis* (D64001). Others include: a sequence from *Synechococcus* (S04426); a sequence (created as Swissprot sequence VIRT18179) in the upstream region of the aspartate racemase gene (D84067) from the archaeon *Desulfurococcus* (Yohda et al. 1996); GDP-mannose pyrophosphorylase (Q46859) from *Yersinia enterocolitica*; and P07874—the 56-kDa bifunctional enzyme from *Pseudomonas aeruginosa* with both PMI and GDP-mannose pyrophosphorylase activities (May et al. 1994). This latter

enzyme is involved in the polymerization of alginate—a compound that protects the cell from host immune responses and antibiotics and is also the major cause of mortality in patients suffering from cystic fibrosis.

There are also eukaryotic equivalents of these enzymes (e.g., S53039 from yeast) in which either one or both motifs are conserved. Related sequences from *Caenorhabditis elegans* include two with unknown function, gi|1707132 and gi|1082118.

Epimerases Involved in Cell Wall Synthesis

In addition to those of the previous section, another group of capsule enzymes share the two-motif structure. These are the epimerases such as TDP-deoxyglucose epimerase (L09188) from *Neisseria meningitidis* and a similar sequence (L33181) from *Yersinia pseudotuberculosis* (Thorson et al. 1994). A smaller member of this family is the sequence (Z81360) recently reported from *Mycobacterium tuberculosis*. This study also identified one sequence, of animal origin, with particular similarity to the *Neisseria* epimerase described above. This is a DTDP-4-rhamnose-3,5-epimerase (U39645) from *C. elegans*.

Eukaryotic Dioxygenases

One class of these iron-containing enzymes contains the two motifs identified here, although the first motif has only a single histidine. These enzymes are the 3-hydroxyanthranilate-3,4-dioxygenases (3-HAO) that catalyze the synthesis of quinolinic acid from 3-hydroxyanthranilic acid. They include sequences from yeast (P47096), *C. elegans* (Z70755), and rat (D44494) (Malherbe et al. 1994). Additionally, the cysteine dioxygenases such as that (P21816) from rat (Hosokawa et al. 1990) clearly contain an equivalent first motif (Fig. 1), although the second is either absent or present as a weakly similar motif separated by a gap of approximately 40 residues.

Proteins Related to DNA Structure and Metabolism

Use of the consensus of domain 1428 in a BLASTP search shows other nonenzymatic proteins which contain the two motifs. One such example is the sequence (X52890) encoding the 280-AA CelD operon repressor from *E. coli* (Parker and Hall 1990). This protein is a unique member of the ARAC/XYLS family of transcriptional regulators (Bustos and Schleif 1993); exceptionally, it is a repressor protein, whereas all other members are positive regulators. With the exception of the recently identified sequences gi|1781124 from *M. tuberculosis* and gi|1788802 from *E. coli*, the other members of the family contain only the second motif. Three other DNA-binding proteins known to contain only the second motif are zinc-finger proteins from the rat (X59993) (Hoog et

al. 1991), mouse (Z32675), bean (X95995), and *Arabidopsis* (gi|1922960). They all contain this motif about 70 residues from the C-terminus.

Also in this category is an *Arabidopsis* protein (DRT 102) concerned with DNA damage repair/toleration and identified in a search for genes to complement *E. coli* mutants lacking defence against UV-light damage to DNA (Pang et al. 1993). Such a function may be relevant to the discussion on the SASP proteins below.

SASPs and Other Small Proteins

Additional analysis identified a series of progressively smaller proteins, each of which contains motifs 1 and 2. Amongst these proteins, none of which has any known function, are two similar accessions (D90909, D90910) from *Synechocystis*, the first of which has the *Arabidopsis* GLP X91957 as the closest relative, and the second of which is similar to cytochrome c551 from *Rhodococcus*. Slightly shorter sequences include U67602 from predicted coding region MJ1618 of *Methanococcus jannaschii* (Bult et al. 1996) and two *E. coli* sequences, (P38522) similar to an aldehyde dehydrogenase (Heim and Strehler 1991), and its longer version—the immunity repressor protein (D90768). The smallest (77 AA) of all those identified is a membrane-spanning protein (MSP) (Li and Strohl 1996) (U37580) from *Streptomyces coelicolor*. This protein, which has similarity to *E. coli* PMI, and also to an upstream sequence (SwissProt virtual sequence VIRT13653) in the prisma gene (Z11975) from *Desulfovibrio desulfuricans* (Stokkermans et al. 1992), is the putative progenitor of all these two motif proteins.

Subsequently, a BLASTP search revealed similarity (11/29 residues identical) between a region spanning the central (intermotif) portion of the MSP sequence and part of the 71-AA small, acid-soluble spore protein (SASP) from *Thermactinomyces thalophilus* (M13042). This is one member of a family of proteins that bind to spore DNA (double stranded) and cause the DNA to change to an A-like conformation. They thus protect the DNA backbone from enzymic and nonenzymic cleavage (Setlow 1995). Of particular relevance is their resistance to hydrogen peroxide (Popham et al. 1995)—a most unusual feature also found in the cereal oxalate oxidase described above.

Discussion

There are two main conclusions to be drawn from this study. First, it has identified for the first time a series of prokaryotic and eukaryotic microbial proteins (Fig. 2) with strong similarity to the two-domain structure of the storage proteins of higher plants. Whether any oxalate decarboxylase (cf. I25120) or other enzyme activity has been retained during this evolution is unknown. Second,

it provides evidence that amongst a broad range of mostly microbial organisms, there has been evolution, presumably by recombination and minor duplication, to produce several progressively larger proteins, ranging in length from the 77-AA *Streptomyces* membrane-spanning protein to the 1,214-AA rat zinc-finger protein. Each protein contains a conserved sequence, usually including two histidine-containing motifs separated by 15–27 residues (summarized in Table 1); the exact distance between motifs is a diagnostic feature of the specific class of protein. It is known that homologous proteins can evolve either different (Murzin 1993) or related (Babbitt et al. 1995) enzymatic activities, but this series is possibly unique in the range of function ascribed to the various proteins—namely, types of isomerase, epimerase, cyclase, oxidase, dioxygenase, decarboxylase, and dehydrogenase, as well as binding proteins for auxin, sucrose, and DNA. Many of these enzymes are involved in production of the extracellular matrix, particularly as part of the sporulation/encystment process. This observation is in agreement with the results of a much more limited study on *E. coli* (Labedan and Riley 1995) in which most of the sequence-related proteins were also related in cellular function; they concluded that 971 paralogous genes could have been derived from only 204 ancestral genes.

Preliminary predictions of tertiary structure of the smaller single-domain proteins described in this study reveal a predominantly β -strand form (data not shown). Presumably, the terminal α -helices found in the larger proteins (>ca. 150 AA), along with the extended gap (>15 AA) between motifs, were added to the core sequence at a very early stage of evolution. These additions to the ends and to the center of the sequence occurred prior to the development of land plants and were part of the increasing complexity that led eventually to the trimeric structure of the desiccation tolerant storage proteins (approx. molecular mass 150–200 kDa) found in seeds. It can be anticipated (see Gane et al. 1997) that structural studies will soon help to describe the exact role of the two motifs, determine the components which provide the resistance to environmental extremes, and identify the various catalytic residues.

Acknowledgments. J.M. Dunwell should like to thank the BBSRC and ZENCA plc. for financial support during the course of this study.

References

- Babbitt PC, Mrachko GT, Hasson MS, Huisman GW, Kolter R, Ringe D, Petsko GA, Kenyon GL, Gerlt JA (1995) A functionally diverse enzyme superfamily that abstracts the α protons of carboxylic acids. *Science* 267:1159–1161
- Bailey TL, Elkan C (1994) Fitting a mixture model by expectation maximization to discover motifs in biopolymers. *Proc 2nd Int Conf Intell Systems for Mol Biol*. AAAI Press, Menlo Park, California, pp 28–36
- Bäumlein H, Braun H, Kakhovskaya IA, Shutov AD (1995) Seed stor-

- age proteins of spermatophytes share a common ancestor with desiccation proteins of fungi. *J Mol Evol* 41:1070–1075
- Bergh S, Uhlen M (1992) Analysis of a polyketide synthesis-encoding gene cluster of *Streptomyces curacoi*. *Gene* 117:131–136
- Berna A, Bernier F (1997) Regulated expression of a wheat germin gene in tobacco: oxalate oxidase activity and apoplasmic localization of the heterologous protein. *Plant Mol Biol* 33:417–429
- Bernier F, Lemieux G, Pallotta D (1987) Gene families encode the major encystment-specific proteins of *Physarum polycephalum* plasmodia. *Gene* 59:265–277
- Blanco G, Brian P, Pereda A, Méndez C, Salas JA, Chater KF (1993) Hybridization and DNA sequence analyses suggest an early evolutionary divergence of related biosynthetic gene sets encoding polyketide antibiotics and spore pigments in *Streptomyces* spp. *Gene* 130:107–116
- Braun H, Czihal A, Shotov AD, Baumlein H (1996) A vicilin-like seed protein of cycads: similarity to sucrose-binding proteins. *Plant Mol Biol* 31:35–44
- Brown JC, Jones AM (1994) Mapping the auxin-binding site of auxin-binding protein 1. *J Biol Chem* 269:21136–21140
- Bult CJ et al. (1996) Complete genome sequence of the methanogenic archaeon *Methanococcus jannaschii*. *Science* 273:1058–1073
- Bustos SA, Schleif RF (1993) Functional domains of the AraC protein. *Proc Natl Acad Sci* 90:5638–5642
- Domon J-M, Dumas B, Lainé E, Meyer Y, David A, David H (1995) Three glycosylated polypeptides secreted by several embryonic cell lines of pine show highly specific serological affinity to antibodies directed against the wheat germin apoprotein monomer. *Plant Physiol* 108:141–148
- Gane PJ, Dunwell JM, Warwicker J (1997) Modelling based on the structure of vicilins predicts a histidine cluster in the active site of oxalate oxidase. *J Mol Evol* (in press)
- Heim R, Strehler EE (1991) Cloning an *Escherichia coli* gene encoding a protein remarkably similar to mammalian aldehyde dehydrogenases. *Gene* 99:15–23
- Heintzen C, Fischer R, Melzer S, Kappeler S, Apel K, Staiger D (1994) Circadian oscillations of a transcript encoding a germin-like protein that is associated with cell walls in young leaves of the long-day plant *Sinapis alba* L. *Plant Physiol* 106:905–915
- Henikoff S, Henikoff JG (1991) Automated assembly of protein blocks for database searching. *Nucleic Acids Res* 19:6565–6572
- Hoog C, Schalling M, Grunder-Brundell E, Daneholt B (1991) Analysis of a murine germ cell-specific transcript encodes a putative zinc finger protein. *Mol Reprod Dev* 30:173–181
- Hosokawa Y, Matsumoto A, Oka J, Itakura H, Yamaguchi K (1990) Isolation and characterization of a cDNA for rat liver cysteine dioxygenase. *Biochem Biophys Res Commun* 168:473–478
- Kaneko T et al. (1996) Sequence analysis of the genome of the unicellular cyanobacterium *Synechocystis* sp. strain PCC6803. II. Sequence determination of the entire genome and assignment of potential protein-coding regions. *DNA Res* 3:109–136
- Labadan B, Riley M (1995) Gene products of *Escherichia coli*: sequence comparisons and common ancestries. *Mol Biol Evol* 12:980–9878
- Lane BG, Bernier F, Dratewka-Kos E, Shafei R, Kennedy TD, Pyne C, Munro JR, Vaughan T, Walters D, Altomare F (1991) Homologies between members of the germin gene family in hexaploid wheat and similarities between these wheat germins and certain *Physarum* spherulins. *J Biol Chem* 266:10461–10468
- Lane BG, Cuming AC, Fregeau J, Carpita NC, Hurkman WJ, Bernier F, Dratewka-Kos E, Kennedy TD (1992) Germin isoforms are discrete temporal markers of early plant development. *Eur J Biochem* 209:961–969
- Lane BG, Dunwell JM, Ray JA, Schmitt MR, Cuming AC (1993) Germin, a protein marker of early plant development, is an oxalate oxidase. *J Biol Chem* 268:12239–12242
- Lawrence MC, Izard T, Beuchat M, Blagrove RJ, Colman PM (1994) Structure of phaseolin at 2.2 Å resolution: Implications for a common vicilin/legumin structure and the genetic engineering of seed storage proteins. *J Mol Biol* 238:748–776
- Li Y, Strohl WR (1996) Cloning, purification, and properties of a phosphotyrosine protein from *Streptomyces coelicolor* A3(2). *J Bacteriol* 178:136–142
- Malherbe P, Kohler C, Da Prada M, Lang G, Kiefer V, Schwarcz R, Lahm HW, Cesura AM (1994) Molecular cloning and functional expression of human 3-hydroxyanthranilic-acid dioxygenase. *J Biol Chem* 269:13792–13797
- May TB, Shinabarger D, Boyd A, Chakrabarty AM (1994) Identification of amino acid residues involved in the activity of phosphomannose isomerase-guanosine 5'-diphospho-D-mannose pyrophosphorylase. A bifunctional enzyme in the alginate biosynthetic pathway of *Pseudomonas aeruginosa*. *J Biol Chem* 269:4872–4877
- Mehta A, Datta A (1991) Oxalate decarboxylase from *Collybia velutipes*. Purification, characterisation and cloning. *J Biol Chem* 266:23548–23553
- Murzina AG (1993) Can homologous proteins evolve different enzymatic activities? *Trends Biochem Sci* 18:403–405
- Ono M, Sage-Ono K, Inoue M, Kamada H, Harada H (1996) Transient increase in the level of mRNA for a germin-like protein in leaves of the short-day plant *Pharbitis nil* during the photoperiodic induction of flowering. *Plant Cell Physiol* 37:855–861
- Pang Q, Hays JB, Rajagopal I, Schaefer TS (1993) Selection of *Escherichia coli* DNA-damage-sensitive mutants and analysis of two plant cDNAs that appear to express UV-specific dark repair activities. *Plant Mol Biol* 22:411–426
- Parker LL, Hall BG (1990) Characterisation and nucleotide sequence of the cryptic *cel* operon of *Escherichia coli* K12. *Genetics* 124:455–471
- Popham DL, Sengupta S, Setlow P (1995) Heat, hydrogen peroxide, and UV resistance of *Bacillus subtilis* spores with increased core water content and with or without major DNA-binding proteins. *Appl Environ Microbiol* 61:3633–3638
- Proudfoot AEI, Turcatti G, Wells TNC, Payton MA, Smith DJ (1994) Purification, cDNA cloning and heterologous expression of human phosphomannose isomerase. *Eur J Biochem* 219:415–423
- Sau S, Lee CY (1996) Cloning of type 8 capsule genes and analysis of gene clusters of different capsular polysaccharides in *Staphylococcus aureus*. *J Bacteriol* 178:2118–2126
- Setlow P (1995) Mechanisms for the prevention of damage to DNA in spores of *Bacillus* species. *Ann Rev Microbiol* 49:29–54
- Smith CA, Pinkney M, Guiney DG, Thomas CM (1993) The ancestral IncP replication system consisted of contiguous *oriV* and *trfA* segments as deduced from a comparison of the nucleotide sequences of diverse IncP plasmids. *J Gen Microbiol* 139:1761–1766
- Stokkermans JPWG, Pierik AJ, Wolbert RBG, Hagen WR, Van Dongen WMAM, Veeger C (1992) The primary structure of a protein containing a putative [6Fe-6S] prismatic cluster from *Desulfovibrio vulgaris* (Hildenborough). *Eur J Biochem* 208:435–442
- Thorson JS, Lo SF, Ploux O, He X, Liu HW (1994) Studies on the biosynthesis of 3,6-dideoxyhexoses: molecular cloning and characterization of the asc (ascarylose) region from *Yersinia pseudotuberculosis* serogroup VA. *J Bacteriol* 176:5483–5493
- Venis MA, Napier R (1995) Auxin receptors and auxin binding proteins. *Crit Rev Plant Sci* 14:27–47
- Worley KC, Wiese BA, Smith RF (1990) BEAUTY: an enhanced BLAST-based search tool that integrates multiple biological information resources into sequence similarity search results. *Genome Res* 5:173–184
- Yohda M, Endo I, Abe Y, Ohta T, Iida T, Maruyama T, Kagawa Y (1996) Gene for aspartate racemase from the sulfur-dependent hyperthermophilic archaeum, *Desulfurococcus* strain SY. *J Biol Chem* 271:22017–22021
- Yoshida K, Seki S, Fujimura M, Miwa Y, Fujita Y (1995) Cloning and sequencing of a 36-kb region of the *Bacillus subtilis* genome between the *gnt* and *iol* operons. *DNA Res* 2:61–69