

# BIOCHEMICAL JOURNAL LETTERS

## A classification of nucleotide-diphospho-sugar glycosyltransferases based on amino acid sequence similarities

Glycosyl-transfer reactions are, on quantitative terms, the most important biotransformations on Earth, since they account for the biosynthesis and hydrolysis of the bulk of biomass [1]. The biosynthesis of polysaccharides and complex carbohydrates is also of fundamental biological importance, since these molecules of fascinating diversity directly mediate a wide range of functions, from structure and storage to specific signalling. The biosynthesis of disaccharides, oligosaccharides and polysaccharides involves the action of hundreds of different glycosyltransferases (EC 2.4.x.y), enzymes which catalyse the transfer of sugar moieties from activated donor molecules to specific acceptor molecules, forming glycosidic bonds. There is a parallel extensive diversity of glycoside hydrolases (EC 3.2.1.x), enzymes which cleave such bonds to yield carbohydrates smaller than those from whence they originated. The immense functional and structural variety of glycosyltransferases and glycoside hydrolases raises the problem of their classification.

Regardless of the direction of the reaction, enzymes which catalyse glycosyl-transfer reactions can be classified according to the stereochemistries of the reaction substrates and products as either 'retaining' or 'invertin'g' enzymes [2]. Furthermore, specific enzymes can be classified on the basis of the reaction catalysed and the substrate specificity, according to the recommendations of the International Union of Biochemistry and Molecular Biology (IUBMB) [3]. However, there are limitations to the utility of this system for classification of glycosyltransferases and glycoside hydrolases, as it does not indicate the intrinsic structural features of the enzymes, nor does it adequately accommodate enzymes which act on several distinct substrates.

Classification of enzymes based on the similarities of their amino acid sequences offers a system complementary to that of the IUBMB [3] and realizes the potential to marry structural features of enzymes with their observed functions. Such classification systems have been proposed for glycoside hydrolases [4] and peptidases [5], and have been updated with the increasing number of cloned genes for these enzymes [6–9]. A significant advantage of classification according to sequence similarities is that it allows logical grouping of enzymes of different EC numbers into polyspecific families and offers insights into the divergent evolution of enzyme families [4]. Conversely, some enzymes which can be grouped by function have been shown to belong to several distinct families and thus reflect convergent evolution [4]. Significantly, the discriminatory power of these classifications has been confirmed by the similarity of the three-dimensional structures [10] and the conserved molecular mechanisms [11] of family members.

Despite the utility of the sequence-based classification of glycoside hydrolases, no such system has been fully described for glycosyltransferases. One difficulty with such a classification is the number of enzymes concerned (193 entries of EC 2.4.1.x), and the diversity of sugar donors. These can be di- or poly-

saccharides, sugar 1-phosphates, or, most commonly, nucleotide diphospho-sugars (NDP-sugars). Whereas a number of the latter type of glycosyltransferases have been compared and grouped into a single family [12], and sequence similarity has been used to predict mechanisms of action [13], there have been no reports of a comprehensive classification of NDP-sugar glycosyltransferases. The present letter describes a classification of NDP-sugar hexosyltransferases (EC 2.4.1.x) and related proteins into distinct sequence-based families.

Sequences of NDP-sugar hexosyltransferases were retrieved from the SwissProt and EMBL/GenBank databanks and compiled into a preliminary sequence library which covered the 35 EC 2.4.1.x entries for which at least one sequence is known to date. Representatives of each EC number were used as templates for BLAST similarity searches [14], and complementary sequences were retrieved from either SwissProt or EMBL/GenBank. BLAST results were examined using Visual BLAST [15]. When the BLAST probability values were low (typically  $P > 10^{-3}$ ), sequences were further compared by hydrophobic cluster analysis (HCA) [16,17]. A family was defined as a grouping of at least two sequences of significant amino acid or HCA similarity over a length exceeding 100 residues, with no similarity to other families.

A total of 555 sequences were analysed, of which 553 were classified into 26 families (Table 1). Only two sequences, namely those of the mannosyltransferase OCH1 of *Saccharomyces cerevisiae* (GenBank D11095) and the DNA  $\beta$ -glucosyltransferase of bacteriophage T4 (Swiss-Prot P04547), could not be assigned to any family and were left unclassified.

Seven families were found polyspecific (containing two or more EC numbers), whereas the others were either monospecific (one single EC number) or 'uncertain' (no EC numbers assigned to the sequences). Previous experience with the classification of glycoside hydrolases suggests that the number of polyspecific families could increase with the availability of more glycosyltransferase sequences.

More than half of the sequences are found in the three largest families (families 1, 2 and 4 with respectively 107, 139 and 84 members; Table 1). Family 1 comprises proteins from viruses, bacteria, fungi, plants and animals. Families 2, 4, 8 and 20 contain sequences from bacteria, fungi, plants and animals. Conversely, several other small families appear strongly biased toward only one taxonomic group, but this could simply be a consequence of the smaller number of current members in these families.

Sequence similarity is strongly indicative of folding similarity in proteins [18]. Conservation of tertiary structure is such that the same three-dimensional fold is expected to be found within each of the families defined by the present study. For polyspecific families, this suggests that details of the three-dimensional structure, rather than differences in the global fold, will explain different donor and/or acceptor specificities. Whereas to date there has been only one reported three-dimensional structure for a glycosyltransferase, the DNA  $\beta$ -glucosyltransferase of bacteriophage T4 [19], it is inevitable that more of these enzymes will be purified, crystallized and characterized. Family allocation

**Table 1 Families of NDP-sugar hexosyltransferases and related protein sequences**

The mechanism (retaining or inverting) is indicated for each family where it could be unambiguously identified from the sequence databanks or from the EC recommendations. Notes: <sup>(a)</sup>Description as found in the sequence databanks; <sup>(b)</sup>EC number as found in the sequence databanks; <sup>(c)</sup>accession numbers starting with P or Q are from the Swiss-Prot databank, those starting with PC are from the PIR databank and those starting with other letters are from the EMBL/GenBank databanks. For conciseness, the names of several enzymes have been abbreviated when a known EC number could be given. Specifically, this has been done for EC 2.4.1.101 ( $\beta$ -1,3-mannosyl-glycoprotein  $\beta$ -1,2-*N*-acetylglucosaminyltransferase), EC 2.4.1.102 ( $\beta$ -1,3-galactosyl-*O*-glycosyl-glycoprotein  $\beta$ -1,6-*N*-acetylglucosaminyltransferase), EC 2.4.1.143 ( $\alpha$ -1,6-mannosyl-glycoprotein  $\beta$ -1,2-*N*-acetylglucosaminyltransferase), EC 2.4.1.144 ( $\beta$ -1,4-mannosyl-glycoprotein  $\beta$ -1,4-*N*-acetylglucosaminyltransferase), and EC 2.4.1.155 [ $\alpha$ -1,3(6)-mannosyl-glycoprotein  $\beta$ -1,6-*N*-acetylglucosaminyltransferase]. Abbreviations used: GPI, glycosyl-phosphatidylinositol; LPS, lipopolysaccharide; GlcNAc, *N*-acetylglucosamine.

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Family 1 (inverting)			
Glycosyltransferase GtA		<i>Amycolatopsis orientalis</i>	U84349
Glycosyltransferase GtB		<i>Amycolatopsis orientalis</i>	U84349
Glycosyltransferase GtC		<i>Amycolatopsis orientalis</i>	U84349
Glycosyltransferase GtD		<i>Amycolatopsis orientalis</i>	U84350
Glycosyltransferase GtE		<i>Amycolatopsis orientalis</i>	U84350
Ecdysteroid glucosyltransferase		<i>Autographa californica</i> NP virus	P18569
Unknown AC3.2		<i>Caenorhabditis elegans</i>	Z71177
Unknown AC3.7		<i>Caenorhabditis elegans</i>	Z71177
Unknown AC3.8		<i>Caenorhabditis elegans</i>	Z71177
Unknown B0310.5		<i>Caenorhabditis elegans</i>	U40959
Unknown C07A9.6		<i>Caenorhabditis elegans</i>	P34317
Unknown C08B6.1		<i>Caenorhabditis elegans</i>	Z72502
Unknown C17G1.3		<i>Caenorhabditis elegans</i>	Z78415
Unknown C18C4.3		<i>Caenorhabditis elegans</i>	U55369
Unknown C23G10.6		<i>Caenorhabditis elegans</i>	U39851
Unknown C33A12.6		<i>Caenorhabditis elegans</i>	Z68493
Unknown C35A5.2		<i>Caenorhabditis elegans</i>	Z71185
Unknown C44H9.1		<i>Caenorhabditis elegans</i>	Z75529
Unknown C55H1.1		<i>Caenorhabditis elegans</i>	U55367
Unknown F01E11.1		<i>Caenorhabditis elegans</i>	U42832
Unknown F08G5.5		<i>Caenorhabditis elegans</i>	Z70682
Unknown F29F11.2		<i>Caenorhabditis elegans</i>	Z73905
Unknown F35H8.6		<i>Caenorhabditis elegans</i>	Z36752
Unknown R04B5.9		<i>Caenorhabditis elegans</i>	Z70782
Unknown R11A8.3		<i>Caenorhabditis elegans</i>	Z70310
Unknown T04H1.7		<i>Caenorhabditis elegans</i>	Z78200
Unknown T04H1.8		<i>Caenorhabditis elegans</i>	Z78200
Unknown T07C5.1		<i>Caenorhabditis elegans</i>	Z50006
Unknown T25B9.7		<i>Caenorhabditis elegans</i>	Z70311
Unknown ZC443.6		<i>Caenorhabditis elegans</i>	Z75553
Unknown ZC455.3		<i>Caenorhabditis elegans</i>	Z75554
Unknown ZC455.4		<i>Caenorhabditis elegans</i>	Z75554
Unknown ZC455.5		<i>Caenorhabditis elegans</i>	Z75554
Unknown ZC455.6		<i>Caenorhabditis elegans</i>	Z75554
Ecdysteroid glucosyltransferase		<i>Choristoneura fumiferana</i> NP virus 1	U10441
Ecdysteroid glucosyltransferase		<i>Choristoneura fumiferana</i> NP virus 2	U10476
CrtX protein		<i>Erwinia herbicola</i>	M90698
Zeaxanthin glucosyltransferase		<i>Erwinia herbicola</i>	Q01330
Zeaxanthin glucosyltransferase		<i>Erwinia uredovora</i>	P21686
Flavonol <i>O</i> <sup>3</sup> -glucosyltransferase	2.4.1.91	<i>Gentiana triflora</i>	D85186
Flavonol <i>O</i> <sup>3</sup> -glucosyltransferase	2.4.1.91	<i>Hordeum vulgare</i>	P14726
1- $\beta$ -Galactosyltransferase	2.4.1.45	Human	U32370
Glucuronosyltransferase 1B	2.4.1.17	Human	P36509
Glucuronosyltransferase 1C	2.4.1.17	Human	P35503
Glucuronosyltransferase 1D	2.4.1.17	Human	P22310
Glucuronosyltransferase 1E	2.4.1.17	Human	P35504
Glucuronosyltransferase 1A	2.4.1.17	Human	P22309
Glucuronosyltransferase 1F	2.4.1.17	Human	P19224
Glucuronosyltransferase 2B10	2.4.1.17	Human	P36537
Glucuronosyltransferase 2B11	2.4.1.17	Human	P36538
Glucuronosyltransferase 2B4	2.4.1.17	Human	P06133
Glucuronosyltransferase 2B7	2.4.1.17	Human	P16662
Glucuronosyltransferase 2B8	2.4.1.17	Human	P23765
Glucuronosyltransferase 2B15	2.4.1.17	Human	P54855
Ecdysteroid glucosyltransferase		<i>Lacanobia oleracea</i> granulosis virus	Y08294
Twil1 protein		<i>Lycopersicon esculentum</i>	X85138
Ecdysteroid glucosyltransferase		<i>Lymantria dispar</i> multicapsid NP virus	P41713
Ecdysteroid glucosyltransferase		<i>Mamestra brassicae</i> NP virus	U41999
UDP-glucose glucosyltransferase Cgt1		<i>Manihot esculenta</i>	X77459
UDP-glucose glucosyltransferase Cgt5		<i>Manihot esculenta</i>	X77462
Galactosyltransferase	2.4.1.62	Mouse	X92122
Glucuronosyltransferase 1A1	2.4.1.17	Mouse	D87867

Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Glucuronosyltransferase 1-06	2.4.1.17	Mouse	U16818
Glucuronosyltransferase UGTP4	2.4.1.17	Mouse	L27122
Glucuronosyltransferase UGTBr1	2.4.1.17	Mouse	S64760
Glucuronosyltransferase 2B5	2.4.1.17	Mouse	P17717
Unknown (PID: E256387)		<i>Mycobacterium tuberculosis</i>	Z77826
Unknown (PID: E256535)		<i>Mycobacterium tuberculosis</i>	Z77826
Salicylate-induced glucosyltransferase IS10a		<i>Nicotiana tabacum</i>	U32643
Salicylate-induced glucosyltransferase IS5a		<i>Nicotiana tabacum</i>	U32644
UDP-rhamnose rhamnosyltransferase		<i>Petunia hybrida</i>	Z25802
Glucuronosyltransferase		<i>Pleuronectes platessa</i>	X74116
Rhamnosyl transferase		<i>Pseudomonas aeruginosa</i>	L28170
Glucuronosyltransferase UGT1-6		Rabbit	U09030
Glucuronosyltransferase UGT1-4		Rabbit	U09101
Glucuronosyltransferase 2A2		Rabbit	P36514
Glucuronosyltransferase 2B13		Rabbit	P36512
Glucuronosyltransferase 2B14		Rabbit	P36513
1- $\beta$ -Galactosyltransferase	2.4.1.45	Rat	Q09426
Glucuronosyltransferase 1.1	2.4.1.17	Rat	U20551
Glucuronosyltransferase		Rat	J05132
Glucuronosyltransferase 2B12		Rat	P36511
Glucuronosyltransferase		Rat	U27518
Glucuronosyltransferase UGT1		Rat	D38063
Glucuronosyltransferase 1A		Rat	P20720
Glucuronosyltransferase 1F		Rat	P08430
Glucuronosyltransferase 2A1		Rat	P36510
Glucuronosyltransferase 2B1		Rat	P09875
Glucuronosyltransferase 2B2		Rat	P08541
Glucuronosyltransferase 2B3		Rat	P08542
Glucuronosyltransferase 2B6		Rat	P19488
Glucuronosyltransferase R-21		Rat	P19489
Unknown (gene L9470.23)		<i>Saccharomyces cerevisiae</i>	U17246
Glycosyltransferase		<i>Solanum melongena</i>	X77369
Glucuronosyltransferase		<i>Solanum tuberosum</i>	U82367
Ecdysteroid glucosyltransferase		<i>Spodoptera littoralis</i> NP virus	X84701
Glycosyltransferase		<i>Streptomyces antibioticus</i>	Z22577
Glycosyltransferase		<i>Streptomyces fradiae</i>	X81885
Glycosyltransferase		<i>Streptomyces peucetius</i>	L47164
Macrolide glycosyltransferase		<i>Streptomyces lividans</i>	M74717
Glycosyltransferase		<i>Streptomyces</i> sp. C5	U43704
Zeaxanthin glucosyltransferase CrtX		<i>Synechocystis</i> sp.	D90899
Flavonol O <sup>2</sup> -glucosyltransferase	2.4.1.91	<i>Vitis vinifera</i>	P51094
Flavonol O <sup>2</sup> -glucosyltransferase 1	2.4.1.91	<i>Zea mays</i>	P16165
Flavonol O <sup>2</sup> -glucosyltransferase 2	2.4.1.91	<i>Zea mays</i>	P16166
Flavonol O <sup>2</sup> -glucosyltransferase 3	2.4.1.91	<i>Zea mays</i>	P16167
Indole-3-acetate $\beta$ -glucosyltransferase	2.4.1.121	<i>Zea mays</i>	L34847
Family 2 (inverting)			
Cellulose synthase AcsA	2.4.1.12	<i>Acetobacter xylinum</i>	P19449
Cellulose synthase AcsAll	2.4.1.12	<i>Acetobacter xylinum</i>	U15957
Cellulose synthase BcsA	2.4.1.12	<i>Acetobacter xylinum</i>	P21877
Cellulose synthase	2.4.1.12	<i>Agrobacterium tumefaciens</i>	L38609
Unknown		<i>Anabaena</i> sp.	P22639
Chitin synthase C	2.4.1.16	<i>Aspergillus fumigatus</i>	X94245
Chitin synthase G-1	2.4.1.16	<i>Aspergillus fumigatus</i>	U39478
Chitin synthase G-2	2.4.1.16	<i>Aspergillus fumigatus</i>	U39479
GlcNAc transferase (NodC)		<i>Azorhizobium caulinodans</i>	Q07755
CgeBB protein		<i>Bacillus subtilis</i>	P42092
CsbB protein		<i>Bacillus subtilis</i>	L77099
Teichoic acid biosynthesis protein GgaA		<i>Bacillus subtilis</i>	P46917
Teichoic acid biosynthesis protein GgaB		<i>Bacillus subtilis</i>	P46918
SpsA protein		<i>Bacillus subtilis</i>	P39621
YwdF protein		<i>Bacillus subtilis</i>	P39614
YveO protein		<i>Bacillus subtilis</i>	Z71928
Polypeptide GalNAc transferase	2.4.1.41	Bovine	Q07537
GlcNAc transferase (NodC)		<i>Bradyrhizobium elkanii</i>	D28963
GlcNAc transferase (NodC)		<i>Bradyrhizobium</i> sp.	P04677
Unknown F13G3.6		<i>Caenorhabditis elegans</i>	Z71259
Unknown ZK688.8		<i>Caenorhabditis elegans</i>	P34678
Chitin synthase	2.4.1.16	<i>Candida albicans</i>	P30572
Chitin synthase	2.4.1.16	<i>Candida albicans</i>	P30573
Chitin synthase	2.4.1.16	<i>Candida maltosa</i>	D29762

Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Chitin oligosaccharide synthase		<i>Danio rerio</i>	U53223
Chitin synthase A	2.4.1.16	<i>Emericella nidulans</i>	D21268
Chitin synthase B	2.4.1.16	<i>Emericella nidulans</i>	D21269
Chitin synthase D	2.4.1.16	<i>Emericella nidulans</i>	U62895
Chitin synthase E	2.4.1.16	<i>Emericella nidulans</i>	U52362
AmsB protein		<i>Erwinia amylovora</i>	X77921
AmsE protein		<i>Erwinia amylovora</i>	X77921
KfiC protein		<i>Escherichia coli</i>	X77617
WcaA protein		<i>Escherichia coli</i>	U38473
WcaE protein		<i>Escherichia coli</i>	U38473
Dolichol-phosphate mannosyltransferase	2.4.1.83	<i>Escherichia coli</i>	D90856
YhjO protein		<i>Escherichia coli</i>	P37653
Unknown (ORF275)		<i>Escherichia coli</i>	L04596
Unknown protein F344		<i>Escherichia coli</i>	P11290
Cellulose synthase	2.4.1.12	<i>Gossypium hirsutum</i>	U58283
Unknown HI0868		<i>Haemophilus influenzae</i>	U32768
Unknown HI1578		<i>Haemophilus influenzae</i>	U32832
Unknown HI1695		<i>Haemophilus influenzae</i>	U32842
Unknown HI1696		<i>Haemophilus influenzae</i>	U32842
Unknown (ORF5)		<i>Haemophilus influenzae</i>	M94855
Unknown (ORF6)		<i>Haemophilus influenzae</i>	M94855
Hyaluronan synthase 1		Human	D84424
Hyaluronan synthase 2		Human	U59269
Polypeptide GalNAc transferase T1	2.4.1.41	Human	Q10472
Polypeptide GalNAc transferase T2	2.4.1.41	Human	X85019
Polypeptide GalNAc transferase T3	2.4.1.41	Human	X92689
EpsG protein		<i>Lactococcus lactis cremoris</i>	U93364
Dolichol-phosphate mannosyltransferase	2.4.1.83	<i>Methanococcus jannaschii</i>	U67504
Unknown (MJ1057)		<i>Methanococcus jannaschii</i>	U67549
Chitin oligosaccharide synthase		Mouse	U53222
Hyaluronan synthase 1		Mouse	D82964
Hyaluronan synthase 2		Mouse	U52524
Polypeptide GalNAc transferase	2.4.1.41	Mouse	U70538
Unknown (PID:G699128)		<i>Mycobacterium leprae</i>	U15180
Unknown (PID:E264129)		<i>Mycobacterium tuberculosis</i>	Z79701
Unknown (PID:E264145)		<i>Mycobacterium tuberculosis</i>	Z79701
Unknown (PID:E264147)		<i>Mycobacterium tuberculosis</i>	Z79701
Unknown (PID:E256390)		<i>Mycobacterium tuberculosis</i>	Z77826
Unknown (PID:E283381)		<i>Mycobacterium tuberculosis</i>	Z83018
RfbV protein		<i>Mycoplasma genitalium</i>	U39685
TrsB protein		<i>Mycoplasma genitalium</i>	P47271
RfbC protein		<i>Myxococcus xanthus</i>	U36795
GlcNAc transferase LgtA		<i>Neisseria gonorrhoeae</i>	U14554
GalNAc transferase LgtD		<i>Neisseria gonorrhoeae</i>	U14554
Glycosyltransferase LgtA		<i>Neisseria meningitidis</i>	U25839
Chitin synthase 1	2.4.1.16	<i>Neurospora crassa</i>	P30588
Chitin synthase 2	2.4.1.16	<i>Neurospora crassa</i>	P30589
Chitin synthase 3	2.4.1.16	<i>Neurospora crassa</i>	P29070
Chitin synthase 4	2.4.1.16	<i>Neurospora crassa</i>	Q01285
Unknown		<i>Paramecium bursaria</i>	U42580
Polypeptide GalNAc transferase	2.4.1.41	<i>Chlorella virus 1</i>	U42580
Unknown (ORF1)		Pig	D85389
Alginate-synthesis related protein Alg8		<i>Porphyromonas gingivalis</i>	U60208
MigA protein		<i>Pseudomonas aeruginosa</i>	L22611
Polypeptide GalNAc transferase	2.4.1.41	<i>Pseudomonas aeruginosa</i>	U70729
GlcNAc transferase (NodC)		Rat	Q10473
GlcNAc transferase (NodC)		<i>Rhizobium galegae</i>	P50356
GlcNAc transferase (NodC)		<i>Rhizobium leguminosarum (viciae)</i>	P04340
GlcNAc transferase (NodC)		<i>Rhizobium leguminosarum (phaseoli)</i>	P24151
GlcNAc transferase (NodC)		<i>Rhizobium loti</i>	P17862
ExoM protein		<i>Rhizobium meliloti</i>	P33695
ExoO protein		<i>Rhizobium meliloti</i>	P33697
ExoU protein		<i>Rhizobium meliloti</i>	P33700
ExoW protein		<i>Rhizobium meliloti</i>	P33702
GlcNAc transferase (NodC)		<i>Rhizobium meliloti</i>	P04341
GlcNAc transferase (NodC)		<i>Rhizobium sp. strain N33</i>	U53327
GlcNAc transferase (NodC)		<i>Rhizobium sp. strain NGR234</i>	P50357
GlcNAc transferase (NodC)		<i>Rhizobium tropici</i>	X98514
Chitin synthase	2.4.1.16	<i>Rhizopus oligosporus</i>	P30594
Unknown		<i>Rhodococcus sp.</i>	P46370
Chitin synthase 1	2.4.1.16	<i>Saccharomyces cerevisiae</i>	P08004
Chitin synthase 2	2.4.1.16	<i>Saccharomyces cerevisiae</i>	P14180

Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Chitin synthase 3	2.4.1.16	<i>Saccharomyces cerevisiae</i>	P29465
Dolichol-phosphate mannosyltransferase	2.4.1.83	<i>Saccharomyces cerevisiae</i>	P14020
Dolichyl-phosphate $\beta$ -glucosyltransferase	2.4.1.117	<i>Saccharomyces cerevisiae</i>	P40350
Rhamnosyltransferase RfbQ		<i>Salmonella enterica</i>	X61917
Rhamnosyltransferase WbaN		<i>Salmonella enterica</i>	X60665
<i>N</i> -Acetylmannosamine transferase RfbA		<i>Salmonella enterica</i> (plasmid pWQ799)	L39794
Rhamnosyltransferase RfbN		<i>Salmonella typhimurium</i>	P26403
RfbV protein		<i>Salmonella typhimurium</i>	P26401
Unknown (ORF 14.1)		<i>Salmonella typhimurium</i>	M65054
Chitin synthase	2.4.1.16	<i>Saprolegnia monoica</i>	U19946
Chitin synthase (class VI)	2.4.1.16	<i>Sartorya fumigata</i>	U62614
KdtX protein		<i>Serratia marcescens</i>	U52844
RfpA protein		<i>Shigella dysenteriae</i>	S73325
RfbG protein		<i>Shigella flexneri</i>	X71970
Unknown (ORF7)		<i>Shigella sonnei</i>	U34305
SpsL protein		<i>Sphingomonas</i> S88	U51197
IcaA protein		<i>Staphylococcus epidermidis</i>	U43366
Fbf15 protein		<i>Stigmatella aurantiaca</i>	Z11601
Cellobiouronic acid synthase		<i>Streptococcus pneumoniae</i>	Z47210
Type 3 capsular polysaccharide synthase		<i>Streptococcus pneumoniae</i>	U15171
Hyaluronan synthase 1		<i>Streptococcus pyogenes</i>	L20853
Hyaluronan synthase		<i>Streptococcus pyogenes</i> strain WF14	L21187
EpsI protein		<i>Streptococcus thermophilus</i>	U40830
Unknown (ORF C04008)		<i>Sulfolobus solfataricus</i>	Y08257
Unknown		<i>Synechococcus</i> sp.	P42460
Unknown (ORF SLL1664)		<i>Synechocystis</i> sp.	D90900
Unknown (ORF SLL1020)		<i>Synechocystis</i> sp.	D90901
Unknown (ORF SLR1537)		<i>Synechocystis</i> sp.	D90906
SpsA protein		<i>Synechocystis</i> sp.	D90911
Unknown (ORF SLR2120)		<i>Synechocystis</i> sp.	D90911
Unknown (ORF SLL1377)		<i>Synechocystis</i> sp.	D90912
Dolichol-phosphate mannosyltransferase	2.4.1.83	<i>Trypanosoma brucei</i>	Z54162
Dolichol-phosphate mannosyltransferase	2.4.1.83	<i>Ustilago maydis</i>	P54856
VirB protein		<i>Vibrio anguillarum</i>	L08012
Glycosyltransferase B (putative)		<i>Vibrio cholerae</i>	U72485
DG42 protein		<i>Xenopus laevis</i>	P13563
RfbB protein		<i>Yersinia enterocolitica</i>	Z18920
RfbH protein		<i>Yersinia enterocolitica</i>	Z18920
TrsB protein		<i>Yersinia enterocolitica</i>	Z47767
TrsC protein		<i>Yersinia enterocolitica</i>	Z47767
Unknown (ORF7.8)		<i>Yersinia enterocolitica</i>	U46859
Unknown (ORF10.9)		<i>Yersinia enterocolitica</i>	U46859
HmsR protein		<i>Yersinia pestis</i>	U22837
Abequosyltransferase (putative)		<i>Yersinia pseudotuberculosis</i>	L01777
Family 3 (retaining)			
Glycogen synthase	2.4.1.11	Human	P13807
Glycogen synthase	2.4.1.11	Mouse	P54859
Glycogen synthase	2.4.1.11	Rabbit	P13834
Glycogen synthase	2.4.1.11	Rat	P17625
Glycogen synthase isoform 1	2.4.1.11	<i>Saccharomyces cerevisiae</i>	P23337
Glycogen synthase isoform 2	2.4.1.11	<i>Saccharomyces cerevisiae</i>	P27472
Family 4 (retaining)			
AceC protein		<i>Acetobacter xylinum</i>	X94981
Sucrose synthase	2.4.1.13	<i>Alnus glutinosa</i>	P49034
HepB protein		<i>Anabaena</i> sp.	U68035
Sucrose synthase	2.4.1.13	<i>Arabidopsis thaliana</i>	P49040
Teichoic acid biosynthesis protein E	2.4.1.52	<i>Bacillus subtilis</i>	P13484
Unknown		<i>Bacillus subtilis</i>	P46915
Unknown		<i>Bacillus subtilis</i>	P42982
YveP protein		<i>Bacillus subtilis</i>	Z71928
YggM protein		<i>Bacillus subtilis</i>	D84432
Sucrose-phosphate synthase	2.4.1.14	<i>Beta vulgaris</i>	P49031
Sucrose synthase SBSS1	2.4.1.13	<i>Beta vulgaris</i>	X81974
BpIH protein		<i>Bordetella pertussis</i>	X90711
Galactosyltransferase		<i>Campylobacter hyoilei</i>	X91081
AmsD protein		<i>Erwinia amylovora</i>	X77921
AmsK protein		<i>Erwinia amylovora</i>	X77921
LPS 1,6-galactosyltransferase RfaB		<i>Escherichia coli</i>	P27127
LPS core biosynthesis protein RfaG		<i>Escherichia coli</i>	P25740

Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Mannosyltransferase A		<i>Escherichia coli</i>	D43637
Mannosyltransferase B		<i>Escherichia coli</i>	D43637
Mannosyltransferase C		<i>Escherichia coli</i>	D43637
O-antigen biosynthesis protein rfb		<i>Escherichia coli</i>	X59852
WcaC protein		<i>Escherichia coli</i>	U38473
WcaL protein		<i>Escherichia coli</i>	U38473
Unknown (ORF3)		<i>Haemophilus influenzae</i>	U36398
Sucrose synthase 1	2.4.1.13	<i>Hordeum vulgare</i>	P31922
Sucrose synthase 2	2.4.1.13	<i>Hordeum vulgare</i>	P31923
GlcNAc-phosphatidylinositol transferase		Human	P37287
Galactosyltransferase RfbF		<i>Klebsiella pneumoniae</i>	L41518
Unknown (ORF7)		<i>Klebsiella pneumoniae</i>	D21242
Unknown (ORF8)		<i>Klebsiella pneumoniae</i>	D21242
Capsular LPS biosynthesis protein M		<i>Methanococcus jannaschii</i>	U67549
Unknown (MJ1185)		<i>Methanococcus jannaschii</i>	U67559
LPS biosynthesis related RfbU-protein		<i>Methanococcus jannaschii</i>	U67601
GPI anchor biosynthesis protein PigA		Mouse	D26047
Unknown (U2168F)		<i>Mycobacterium leprae</i>	P54138
Unknown (MTCY20G9.12)		<i>Mycobacterium tuberculosis</i>	Q11152
Unknown (MTCY25D10.36)		<i>Mycobacterium tuberculosis</i>	Z95558
$\alpha$ -1,2-GlcNAc transferase RfaK		<i>Neisseria meningitidis</i>	U35713
LPS yscLPS glycosyltransferase lcsA		<i>Neisseria meningitidis</i>	U39810
Sucrose synthase 1	2.4.1.13	<i>Oryza sativa</i>	P30298
Sucrose synthase 3	2.4.1.13	<i>Oryza sativa</i>	L03366
Sucrose-phosphate synthase	2.4.1.14	<i>Phaseolus aureus</i>	Q01390
Sucrose synthase	2.4.1.13	<i>Pisum sativum</i>	X98598
CpsF protein		<i>Proteus mirabilis</i>	L36873
GlcNAc-phosphatidylinositol transferase		<i>Saccharomyces cerevisiae</i>	P32363
Unknown (ORF YPL175w)		<i>Saccharomyces cerevisiae</i>	Z73531
Mannosyltransferase WbaW		<i>Salmonella enterica</i>	X61917
Mannosyltransferase WbaZ		<i>Salmonella enterica</i>	X61917
LPS GlcNAc transferase RfaK	2.4.1.56	<i>Salmonella typhimurium</i>	P26470
Mannosyltransferase RfbU		<i>Salmonella typhimurium</i>	P26402
Polysaccharide biosynthesis protein VipC		<i>Salmonella typhimurium</i>	Q04975
Galactosyltransferase RfbF		<i>Serratia marcescens</i>	L34167
Galactosyltransferase RfpB		<i>Shigella dysenteriae</i>	S73325
ExpA3 protein		<i>Sinorhizobium meliloti</i>	Z79692
Sucrose-phosphate synthase	2.4.1.14	<i>Solanum tuberosum</i>	X73477
Sucrose synthase	2.4.1.13	<i>Solanum tuberosum</i>	P10691
Sucrose-phosphate synthase	2.4.1.14	<i>Spinacia oleracea</i>	P31928
CapM protein		<i>Staphylococcus aureus</i>	P39862
Cap1G protein		<i>Streptococcus pneumoniae</i>	Z83335
EpsF protein		<i>Streptococcus thermophilus</i>	U40830
EpsG protein		<i>Streptococcus thermophilus</i>	U40830
LPS glycosyltransferase lcsA		<i>Synechocystis</i> sp.	D90906
Mannosyltransferase RfbW		<i>Synechocystis</i> sp.	D64000
Mannosyltransferase RfbU		<i>Synechocystis</i> sp.	D90901
Mannosyltransferase MtfB		<i>Synechocystis</i> sp.	D90911
Unknown (ORF SLR1508)		<i>Synechocystis</i> sp.	D90911
Unknown (ORF SLR0384)		<i>Synechocystis</i> sp.	D63999
Sucrose-phosphate synthase (ORF SLL0045)		<i>Synechocystis</i> sp.	D64006
Unknown (ORF SLR1077)		<i>Synechocystis</i> sp.	D90901
Unknown (ORF SLR1065)		<i>Synechocystis</i> sp.	D90901
Unknown (ORF SLL1971)		<i>Synechocystis</i> sp.	D90905
Unknown (ORF SLL1723)		<i>Synechocystis</i> sp.	D90906
Unknown (ORF SLR1166)		<i>Synechocystis</i> sp.	D90913
Sucrose synthase type 1	2.4.1.13	<i>Triticum aestivum</i>	M26671
Sucrose synthase type 2	2.4.1.13	<i>Triticum aestivum</i>	M26672
Sucrose synthase	2.4.1.13	<i>Tulipa gesneriana</i>	X96939
Sucrose-phosphate synthase	2.4.1.14	<i>Vicia faba</i>	Z56278
GumH protein		<i>Xanthomonas campestris</i>	U22511
RfbC protein		<i>Yersinia enterocolitica</i>	Z18920
RfbP protein		<i>Yersinia enterocolitica</i>	U46859
TrsD protein		<i>Yersinia enterocolitica</i>	Z47767
TrsE protein		<i>Yersinia enterocolitica</i>	Z47767
TrsH protein		<i>Yersinia enterocolitica</i>	Z47767
Sucrose synthase 1	2.4.1.13	<i>Zea mays</i>	P04712
Sucrose synthase 2	2.4.1.13	<i>Zea mays</i>	P49036
Family 5 (retaining)			
Bacterial glycogen synthase	2.4.1.21	<i>Agrobacterium tumefaciens</i>	P39670

Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Bacterial glycogen synthase	2.4.1.21	<i>Bacillus subtilis</i>	P39125
Bacterial glycogen synthase	2.4.1.21	<i>Escherichia coli</i>	P08323
Bacterial glycogen synthase (putative)		<i>Haemophilus influenzae</i>	P45179
Starch synthase	2.4.1.21	<i>Hordeum vulgare</i>	P09842
Starch synthase	2.4.1.21	<i>Ipomoea batatas</i>	U44126
Starch synthase	2.4.1.21	<i>Manihot esculenta</i>	X74160
Bacterial glycogen synthase	2.4.1.21	<i>Methanococcus jannaschii</i>	U67600
Glucosyl transferase		<i>Oryza glaberrima</i>	D10472
Starch synthase (soluble)	2.4.1.21	<i>Oryza sativa</i>	D16202
Starch synthase	2.4.1.21	<i>Oryza sativa</i>	P19395
Starch synthase (granule-bound)	2.4.1.21	<i>Oryza sativa</i>	X64108
Starch synthase (WaxyG)	2.4.1.21	<i>Oryza sativa</i>	X65183
Starch synthase (isoform 1)	2.4.1.21	<i>Pisum sativum</i>	X88790
Starch synthase (isoform 2)	2.4.1.21	<i>Pisum sativum</i>	X88789
Starch synthase (isoform 1)	2.4.1.11	<i>Solanum tuberosum</i>	Q00775
Starch synthase (isoform 2)	2.4.1.21	<i>Solanum tuberosum</i>	X83220
Starch synthase (isoform 3)	2.4.1.21	<i>Solanum tuberosum</i>	X95759
Starch synthase (isoform 4)	2.4.1.21	<i>Solanum tuberosum</i>	X87988
Starch synthase	2.4.1.21	<i>Sorghum bicolor</i>	U23945
Unknown (ORF SLL1393)		<i>Synechocystis</i> sp.	D90899
Unknown (ORF SLL0945)		<i>Synechocystis</i> sp.	D90915
Starch synthase (isoform 1)	2.4.1.11	<i>Triticum aestivum</i>	P27736
Starch synthase (soluble)	2.4.1.21	<i>Triticum aestivum</i>	U48227
Starch synthase (isoform 2)	2.4.1.21	<i>Triticum aestivum</i>	U66377
Starch synthase	2.4.1.21	<i>Zea mays</i>	P04713
Family 6 (retaining)			
Histo-blood-group-A transferase		Baboon	PC1172
Histo-blood-group-B transferase		Baboon	PC1173
3- $\alpha$ -Galactosyltransferase	2.4.1.124	Bovine	P14769
Histo-blood-group-2 transferase		Chimpanzee	PC1166
Histo-blood-group transferase		Crab-eating macaque	PC1171
Histo-blood-group transferase		Dog	U66140
Histo-blood-group-2 transferase		Gorilla	PC1168
Blood-group-B $\alpha$ -1,3-galactosyltransferase		Human	X91874
Fucosylglycoprotein $\alpha$ -GalNAc transferase		Human	P16442
Histo-blood-group-A transferase		Human	X84746
Histo-blood-group-A transferase		Human	J05175
Histo-blood-group-B transferase		Human	PC1165
$\alpha$ -1,3-Galactosyltransferase	2.4.1.151	Marmoset	S71333
$\alpha$ -1,3-Galactosyltransferase		Mouse	M85153
$\alpha$ -1,3-Galactosyltransferase	2.4.1.151	Mouse	P23336
Histo-blood-group-1 transferase		Orang-utan	PC1169
Histo-blood-group-2 transferase		Orang-utan	PC1170
$\alpha$ -1,3-Galactosyltransferase 1		Pig	L36535
$\alpha$ -1,3-Galactosyltransferase 2		Pig	P50127
Family 7 (inverting)			
Glycoprotein 4- $\beta$ -galactosyltransferase	2.4.1.38/90	Bovine	P08037
N-Acetyl-lactosamine synthase		<i>Caenorhabditis elegans</i>	X98132
Unknown R10E11.4		<i>Caenorhabditis elegans</i>	Z29095
Unknown W02B12.11		<i>Caenorhabditis elegans</i>	Z66521
$\beta$ -1,4-Galactosyltransferase CKI		Chicken	U19890
$\beta$ -1,4-Galactosyltransferase CKII		Chicken	U19889
$\beta$ -1,4-Galactosyltransferase		Human	D29805
$\beta$ -1,4-Galactosyltransferase		Human	U10473
$\beta$ -1,4-Galactosyltransferase Gtn2		Human	U10472
$\beta$ -1,4-Galactosyltransferase Gtn6		Human	U10474
Glycoprotein 4- $\beta$ -galactosyltransferase	2.4.1.38/90	Human	P15291
$\beta$ -1,4-Galactosyltransferase (isoform 1)		<i>Lymnaea stagnalis</i>	X99318
$\beta$ -1,4-Galactosyltransferase (isoform 2)		<i>Lymnaea stagnalis</i>	X99319
$\beta$ -1,4-GlcNAc transferase		<i>Lymnaea stagnalis</i>	X80228
Glycoprotein 4- $\beta$ -galactosyltransferase	2.4.1.38	Mouse	P15535
$\beta$ -1,4-Galactosyltransferase		Pig	U63019
Family 8 (retaining)			
Unknown (ORF2)		<i>Bacillus subtilis</i>	P25148
Unknown F56B6.4		<i>Caenorhabditis elegans</i>	U64599
Unknown T10B10.8		<i>Caenorhabditis elegans</i>	Z72514
LPS $\alpha$ -1,2-galactosyltransferase	2.4.1.58	<i>Escherichia coli</i>	P27129
LPS $\alpha$ -1,3-galactosyltransferase	2.4.1.44	<i>Escherichia coli</i>	P27128

**Table 1** (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
$\alpha$ -Galactosyltransferase LgtC		<i>Haemophilus influenzae</i>	P44597
Glycogenin	2.4.1.186	Human	P46976
RfbC protein 1		<i>Klebsiella pneumoniae</i>	L31762
RfbC protein 2		<i>Klebsiella pneumoniae</i>	L41518
Glycosyltransferase LgtC		<i>Neisseria gonorrhoeae</i>	U14554
Glycosyltransferase LgtC		<i>Neisseria meningitidis</i>	U65788
WSI76 water-stress protein		<i>Oryza sativa</i>	D26537
Glycogenin	2.4.1.186	Rabbit	P13280
Galactosyltransferase LgtA		<i>Rhizobium leguminosarum</i>	X94963
Unknown (PID:G152040)		<i>Rhodobacter sphaeroides</i>	M89780
Glycogenin		<i>Saccharomyces cerevisiae</i>	P36143
Unknown 44.5 kDa protein		<i>Saccharomyces cerevisiae</i>	P47011
Unknown YKRO58W		<i>Saccharomyces cerevisiae</i>	Z28283
LPS $\alpha$ -1,2-galactosyltransferase	2.4.1.58	<i>Salmonella typhimurium</i>	P19817
LPS $\alpha$ -1,3-galactosyltransferase	2.4.1.44	<i>Salmonella typhimurium</i>	P19816
Family 9 (inverting)			
RfaC protein		<i>Bordetella pertussis</i>	X90711
LPS 1,2-GlcNAc transferase	2.4.1.56	<i>Escherichia coli</i>	P27242
LPS heptosyltransferase RfaF		<i>Escherichia coli</i>	P37692
LPS heptosyltransferase RfaC		<i>Escherichia coli</i>	P24173
LPS biosynthesis protein LbgB		<i>Haemophilus ducreyi</i>	U58147
OpsX protein		<i>Haemophilus influenzae</i>	U32712
RfaF protein		<i>Haemophilus influenzae</i>	L76100
Unknown HI0523		<i>Haemophilus influenzae</i>	P44011
RfaC protein		<i>Neisseria gonorrhoeae</i>	U10385
RfaC protein 1		<i>Neisseria meningitidis</i>	U40862
RfaC protein 2		<i>Neisseria meningitidis</i>	U35454
RfaC protein		<i>Salmonella typhimurium</i>	P26469
Family 10 (inverting)			
$\alpha$ -(1,3/4)-Fucosyltransferase 3		Bovine	Q11126
Unknown T05A7.5 (tandem repeat)		<i>Caenorhabditis elegans</i>	U40028
Unknown K08F8.3		<i>Caenorhabditis elegans</i>	Z66497
$\alpha$ -1,3-Fucosyltransferase		Chicken	U73678
Galactoside 3(4)-L-fucosyltransferase 3	2.4.1.65	Human	P21217
$\alpha$ -1,3-Fucosyltransferase 4		Human	P22083
$\alpha$ -1,3-Fucosyltransferase 5		Human	Q11128
$\alpha$ -1,3-Fucosyltransferase 6		Human	P51993
$\alpha$ -1,3-Fucosyltransferase 7		Human	Q11130
$\alpha$ -1,3/4-Fucosyltransferase		Human	D89325
$\alpha$ -1,3-Fucosyltransferase 4		Mouse	Q11127
$\alpha$ -1,3-Fucosyltransferase 7		Mouse	Q11131
$\alpha$ -1,3-Fucosyltransferase		Rat	U58860
Family 11 (inverting)			
FUT2 gene product		Bovine	X99620
Unknown C06E1.7		<i>Caenorhabditis elegans</i>	L16559
Unknown F17B5.e		<i>Caenorhabditis elegans</i>	Z81066
Galactoside 2- $\alpha$ -L-fucosyltransferase 1	2.4.1.69	Human	P19526
$\alpha$ -1,2-Fucosyltransferase 2		Human	D82933
$\alpha$ -1,2-Fucosyltransferase Se2		Human	D89327
$\alpha$ -1,2-Fucosyltransferase Sec2		Human	Q10981
Galactoside 2- $\alpha$ -L-fucosyltransferase	2.4.1.69	Mouse	Y09882
$\alpha$ -1,2-Fucosyltransferase		Pig	L50534
FUT2 gene product		Pig	X99621
Galactoside 2- $\alpha$ -L-fucosyltransferase 1	2.4.1.69	Rabbit	Q10979
Galactoside 2- $\alpha$ -L-fucosyltransferase 2	2.4.1.69	Rabbit	Q10983
Galactoside 2- $\alpha$ -L-fucosyltransferase 3	2.4.1.69	Rabbit	X91269
Galactoside 2- $\alpha$ -L-fucosyltransferase 1	2.4.1.69	Rat	Q10980
Galactoside 2- $\alpha$ -L-fucosyltransferase 2	2.4.1.69	Rat	Q10984
Unknown (ORF11.8)		<i>Yersinia enterocolitica</i>	U46859
Family 12 (inverting)			
$\beta$ -1,4-GalNAc transferase	2.4.1.92	Human	Q00973
$\beta$ -1,4-GalNAc transferase 1	2.4.1.92	Mouse	Q09199
$\beta$ -1,4-GalNAc transferase 2	2.4.1.92	Mouse	Q09200
$\beta$ -1,4-GalNAc transferase	2.4.1.92	Rat	Q10468
Family 13 (inverting)			
Unknown B0416.6		<i>Caenorhabditis elegans</i>	U23516
Unknown F48E3.1		<i>Caenorhabditis elegans</i>	U28735



Table 1 (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Unknown M01F1.1		<i>Caenorhabditis elegans</i>	Z46381
GlcNAc transferase I		<i>Cricetulus griseus</i>	U65791
$\beta$ -1,2-GlcNAc transferase I	2.4.1.101	Human	P26572
$\beta$ -1,2-GlcNAc transferase I	2.4.1.101	Mouse	P27808
$\beta$ -1,2-GlcNAc transferase I	2.4.1.101	Rabbit	P27115
$\beta$ -1,2-GlcNAc transferase I	2.4.1.101	Rat	Q09325
Family 14 (inverting)			
$\beta$ -1,6-GlcNAc transferase	2.4.1.102	Bovine	U41320
Unknown T14B4.9		<i>Caenorhabditis elegans</i>	U50191
Unknown F22D6.11		<i>Caenorhabditis elegans</i>	Z71262
Unknown F30A10.4		<i>Caenorhabditis elegans</i>	Z81072
Unknown F44F4.6		<i>Caenorhabditis elegans</i>	Z37092
Unknown R07B7.6		<i>Caenorhabditis elegans</i>	Z75955
Unknown T09E11.a		<i>Caenorhabditis elegans</i>	Z81147
$\beta$ -1,6-GlcNAc transferase C2GNT	2.4.1.102	Human	Q02742
$\beta$ -1,6-GlcNAc transferase IGNT	2.4.1.150	Human	Q06430
$\beta$ -1,6-GlcNAc transferase C2GNT	2.4.1.102	Mouse	Q09324
Enzymic glycosylation-regulating protein		Rat	S79797
Family 15 (retaining)			
Mannosyltransferase Mnt1	2.4.1.131	<i>Candida albicans</i>	X99619
Mannosyltransferase Mnt2	2.4.1.131	<i>Candida albicans</i>	P46592
Glycolipid 2- $\alpha$ -mannosyltransferase	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P27809
Mannosyltransferase Ktr1	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P27810
Mannosyltransferase Ktr2	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P33550
Mannosyltransferase Ktr3	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P38130
Mannosyltransferase Ktr4	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P38131
Mannosyltransferase Ktr5	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P53966
Mannosyltransferase Ktr6	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P54070
Mannosyltransferase Ktr7	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P40504
Mannosyltransferase Yur1	2.4.1.131	<i>Saccharomyces cerevisiae</i>	P26725
Family 16 (inverting)			
$\beta$ -1,2-GlcNAc transferase II	2.4.1.143	Human	Q10469
$\beta$ -1,2-GlcNAc transferase II	2.4.1.143	Rat	Q09326
Family 17 (inverting)			
$\beta$ -1,4-GlcNAc transferase III	2.4.1.144	Human	Q09327
$\beta$ -1,4-GlcNAc transferase III	2.4.1.144	Mouse	Q10470
$\beta$ -1,4-GlcNAc transferase III	2.4.1.144	Rat	Q02527
Family 18 (inverting)			
$\beta$ -GlcNAc transferase		<i>Cricetulus griseus</i>	U62587
$\beta$ -GlcNAc transferase	2.4.1.155	Human	Q09328
$\beta$ -GlcNAc transferase	2.4.1.155	Rat	Q08834
Family 19			
Lipid A disaccharide synthase LpxB	2.4.1.182	<i>Escherichia coli</i>	P10441
Lipid A disaccharide synthase LpxB	2.4.1.182	<i>Haemophilus influenzae</i>	P45011
LpxB protein		<i>Proteus mirabilis</i>	Y09263
LpxB protein		<i>Synechocystis</i> sp.	D64000
Family 20 (retaining)			
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase	2.4.1.15	<i>Arabidopsis thaliana</i>	Y08568
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps1	2.4.1.15	<i>Aspergillus niger</i>	U07184
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase TpsB	2.4.1.15	<i>Aspergillus niger</i>	U63416
Unknown ZK54.2		<i>Caenorhabditis elegans</i>	U58737
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps1	2.4.1.15	<i>Candida albicans</i>	Y07918
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase OtsA	2.4.1.15	<i>Escherichia coli</i>	P31677
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps1	2.4.1.15	<i>Kluyveromyces lactis</i>	Q07158
OtsA protein		<i>Mycobacterium leprae</i>	U15187
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps1	2.4.1.15	<i>Saccharomyces cerevisiae</i>	Q00764
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps2	2.4.1.15	<i>Saccharomyces cerevisiae</i>	P31688
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps3	2.4.1.15	<i>Saccharomyces cerevisiae</i>	P38426
$\alpha$ , $\alpha$ -Trehalose-phosphate synthase Tps1	2.4.1.15	<i>Schizosaccharomyces pombe</i>	P40387
OtsA protein		<i>Synechocystis</i> sp.	D90913
Family 21 (retaining)			
Unknown F20B4.6		<i>Caenorhabditis elegans</i>	U58735
Unknown T06C12.C		<i>Caenorhabditis elegans</i>	Z81116
Unknown YK29C8.5		<i>Caenorhabditis elegans</i>	U53332
Ceramide glucosyltransferase	2.4.1.80	Human	Q16739
Unknown (ORF SLR0813)		<i>Synechocystis</i> sp.	D90911

**Table 1** (cont.)

Description <sup>(a)</sup>	EC number <sup>(b)</sup>	Organism	Accession no. <sup>(c)</sup>
Family 22			
Unknown C14A4.3		<i>Caenorhabditis elegans</i>	Z49909
Mannosyltransferase		Human	D42138
Mannosyltransferase		<i>Saccharomyces cerevisiae</i>	X96417
Family 23 (inverting)			
Fucosyltransferase (NodZ)		<i>Azorhizobium caulinodans</i>	L18897
Fucosyltransferase (NodZ)		<i>Bradyrhizobium japonicum</i>	L22756
<i>N</i> -Acetyl- $\beta$ -D-glucosaminide $\alpha$ -1,6-fucosyltransferase		Pig	D86723
Family 24			
Unknown C12C8.D		<i>Caenorhabditis elegans</i>	Z81467
UDP-glucose glycoprotein glucosyltransferase		<i>Drosophila melanogaster</i>	U20554
Killer toxin-resistance protein Kre5		<i>Saccharomyces cerevisiae</i>	P22023
UDP-glucose glycoprotein glucosyltransferase		<i>Schizosaccharomyces pombe</i>	U38417
Family 25			
Unknown D2045.9		<i>Caenorhabditis elegans</i>	Z35639
LbgA protein		<i>Haemophilus ducreyi</i>	U58147
Lic2B protein		<i>Haemophilus influenzae</i>	U36398
LPS biosynthesis protein		<i>Haemophilus influenzae</i>	X56903
LPS biosynthesis protein		<i>Haemophilus influenzae</i>	L19441
LPS biosynthesis protein Lex-1		<i>Haemophilus influenzae</i>	U32736
Lex2B protein		<i>Haemophilus influenzae</i>	U05670
Galactosyltransferase LgtB		<i>Neisseria gonorrhoeae</i>	U14554
Galactosyltransferase LgtE		<i>Neisseria gonorrhoeae</i>	U14554
Glycosyltransferase LgtB		<i>Neisseria meningitidis</i>	U25839
Glycosyltransferase LgtE		<i>Neisseria meningitidis</i>	U25839
LpsA protein		<i>Pasteurella haemolytica</i>	U15958
Family 26			
AceB protein		<i>Acetobacter xylinum</i>	X94981
UDP- <i>N</i> -acetylmannosaminuronic acid transferase		<i>Escherichia coli</i>	P27836
Teichoic acid biosynthesis protein A		<i>Escherichia coli</i>	P27620
UDP- <i>N</i> -acetylmannosaminuronic acid transferase		<i>Salmonella typhimurium</i>	P37457
Cps19F protein		<i>Streptococcus pneumoniae</i>	U09239
Unknown (ORF SLR1118)		<i>Synechocystis</i> sp.	D90899
Unknown (ORF SLR1271)		<i>Synechocystis</i> sp.	D90913
GumM protein		<i>Xanthomonas campestris</i>	U22511

described by the present work could provide an aid to structural interpretation and, when more structures become available, suggest possible search models for molecular replacement.

For those families where both the NDP-sugar and the linkage formed are known, i.e., all families except families 19, 22, 24, 25 and 26, the classification based on sequence similarity consistently differentiates retaining from inverting enzymes. This is consistent with the conservation of the catalytic machinery of these enzymes within each family. Almost half of the classified sequences have unknown or uncertain functions (Table 1). A fundamental basis for a classification must be that it has predictive power. The present classification allows the prediction of their global function (i.e. NDP-glycosyltransferase) and product stereochemistry (inverted or retained anomeric configuration).

By analogy to glycoside hydrolases, the catalytic machinery of glycosyltransferases is likely to involve Asp and/or Glu residues whose side chains have the appropriate reactivity to act as the general base for acceptor activation or as the nucleophile for the formation of a glycosyl-enzyme intermediate. Site-directed mutagenesis of ribosyltransferases has shown that specific Glu residues are essential for glycosyltransferase activity [20]. For each family of glycosyltransferases, the list of the invariant Asp or Glu residues is therefore likely to contain catalytic residues. In some of the families we describe there are so few such conserved residues that the catalytic machinery is probably directly identi-

fiable. Examples include families 1, 9 and 11 with one invariant Asp, families 5 and 25 with one invariant Glu, families 3 and 4 with two invariant Glu and families 2, 8, 9 and 20 with two invariant Asp residues.

The EC recommendations place all hexosyltransferases in the same subclass (EC 2.4.1.x), regardless of the sugar donor used by the enzymes. There are clear structural, evolutionary and mechanistic similarities between several glycosyltransferases using glycosides as activated sugar donors and glycoside hydrolases. For example, cyclodextrin glucanotransferases (EC 2.4.1.19) and starch branching enzymes (EC 2.4.1.18) are clearly related to a large number of starch-hydrolysing enzymes forming family 13 of glycoside hydrolases [4,6,7,21]. Similarly, endo-xyloglucan transferases (EC 2.4.1.207) display significant similarities to glycoside hydrolase family 16 members [7]. In contrast, we have been unable to detect any sequence similarity between the NDP-sugar glycosyltransferases we have analysed and glycoside hydrolases. This probably reflects particular constraints on the active site of these glycosyltransferases which must accommodate the bulky NDP-moiety.

On several occasions, we observed that enzymes acting on similar substrates with the same mechanism, and classified in different families, displayed intriguing local similarities which could not be extended to the rest of the sequence. This situation, which perhaps reflects the limitations of sequence comparisons at

very high divergence, is reminiscent of the grouping of glycoside hydrolase families into clans where the only sequence similarity is found around the catalytic machinery [22]. An example of such possibly related families are families 3, 4 and 5, which display limited local similarities. Similarly, families 11 and 23 could perhaps be grouped based on the specific instance of the motif VHVRRTD in a family 23 enzyme (porcine *N*-acetyl- $\beta$ -D-glucosaminide  $\alpha$ -1,6-fucosyltransferase) which is almost identical with one of the three highly conserved motifs in family 11 (VHVRRTD motif). Conversely, the proposed grouping of  $\alpha$ -1,3-fucosyltransferases and  $\alpha$ -1,2-fucosyltransferases [23] cannot be confirmed as the corresponding families (10 and 11) do not bear even one conserved residue. Only structural resolution will allow the reliable grouping of families into 'superfamilies' or 'clans'.

That there are several polyspecific families leads to the proposition that the observed differences in substrate specificity probably reflect divergent evolution from an ancestral form of glycosyltransferase. Conversely, we have identified at least one example of an enzyme activity (lipopolysaccharide 1,2-*N*-acetylglucosaminyltransferase; EC 2.4.1.56) which appears in two distinct families (4 and 9), suggesting that this could constitute an example of convergent evolution.

Genome sequencing projects are increasingly delivering large numbers of potential glycosyltransferase sequences and the present classification that brings together structural, mechanistic and sequence-based information is clearly of biocomputing importance. Significantly, a possible function was recently proposed for secreted Fringe-like signalling molecules based on distant sequence similarity with glycosyltransferase sequences [24]. It is our intention to set up an electronic access to this classification similar to that already implemented for the glycosidases [7].

J. C. thanks the Australian Government Department of Industry, Science and Tourism for financial support of his visit to CERMAV. J.C. also thanks the Sugar Research and Development Corporation for partial funding of this research. G.J.D. is a Royal Society University Research Fellow.

James A. CAMPBELL\*, Gideon J. DAVIES†, Vincent BULONE‡§ and Bernard HENRISSAT‡

\*CSIRO Tropical Agriculture, 306 Carmody Road, St. Lucia, Q 4067, Australia, †Department of Chemistry, University of York, Heslington, York YO1 5DD, U.K., and ‡Centre de Recherches sur les Macromolécules Végétales (Affiliated with the Joseph Fourier University), CNRS, BP 53, F-38041 Grenoble Cedex 9, France

§ To whom correspondence should be addressed.

- Kleene, R. and Berger, E. G. (1993) *Biochim. Biophys. Acta* **1154**, 283–325
- Sinnott, M. L. (1990) *Chem. Rev.* **90**, 1171–1202
- International Union of Biochemistry and Molecular Biology (1992) *Enzyme Nomenclature: Recommendations of the Nomenclature Committee of the International Union of Biochemistry and Molecular Biology*, Academic Press, San Diego
- Henrissat, B. (1991) *Biochem. J.* **280**, 309–316
- Rawlings, N. D. and Barrett, A. J. (1993) *Biochem. J.* **290**, 205–218
- Henrissat, B. and Bairoch, A. (1993) *Biochem. J.* **293**, 781–788
- Henrissat, B. and Bairoch, A. (1996) *Biochem. J.* **316**, 695–696
- Rawlings, N. D. and Barrett, A. J. (1994) *Methods Enzymol.* **244**, 1–15
- Rawlings, N. D. and Barrett, A. J. (1995) *Methods Enzymol.* **248**, 105–120
- Davies, G. and Henrissat, B. (1995) *Structure* **3**, 853–859
- Gebler, J., Gilkes, N. R., Claeysens, M., Wilson, D. B., Béguin, P., Wakarchuk, W. W., Kilburn, D. G., Miller, Jr., R. C., Warren, R. A. J. and Withers, S. G. (1992) *J. Biol. Chem.* **267**, 12559–12561
- Geremia, R. A., Petroni, A., Ielpi, L. and Henrissat, B. (1996) *Biochem. J.* **318**, 133–138
- Saxena, I., Brown, Jr., R. M., Fèvre, M., Geremia, R. A. and Henrissat, B. (1995) *J. Bacteriol.* **177**, 1419–1424
- Altschul, S. F., Gish, W., Miller, W., Myers, E. W. and Lipman, D. J. (1990) *J. Mol. Biol.* **215**, 403–410

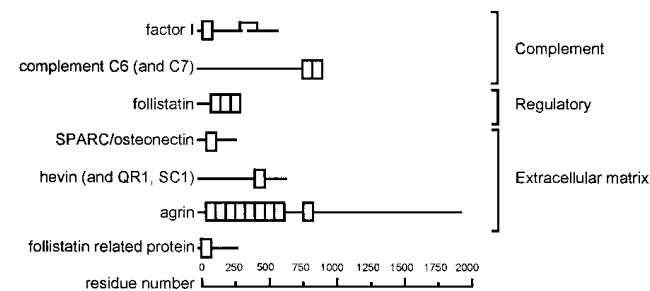
- Durand, P., Canard, L. and Mornon, J. P. (1997) *Comput. Appl. Biosci.*, in the press
- Gaboriaud, C., Bissery, V., Benchetrit, T. and Mornon, J. P. (1987) *FEBS Lett.* **224**, 149–155
- Lemesle-Varloot, L., Henrissat, B., Gaboriaud, C., Bissery, V., Morgat, A. and Mornon, J. P. (1990) *Biochimie* **72**, 555–574
- Chothia, C. and Lesk, A. M. (1986) *EMBO J.* **5**, 823–826
- Vrieland, A., Rüger, W., Driessen, H. P. C. and Freemont, P. S. (1994) *EMBO J.* **13**, 3413–3422
- Takada, T., Iida, K. and Moss, J. (1995) *J. Biol. Chem.* **270**, 541–544
- Jespersen, H. M., MacGregor, E. A., Henrissat, B., Sierks, M. R. and Svensson, B. (1993) *J. Protein Chem.* **12**, 791–805
- Henrissat, B., Callebaut, I., Fabrega, S., Lehn, P., Mornon, J. P. and Davies, G. (1995) *Proc. Natl. Acad. Sci. U.S.A.* **92**, 7090–7094
- Breton, C., Oriol, R. and Imberty, A. (1996) *Glycobiology* **6**, vii–xii
- Yuan, Y. P., Schultz, J., Mlodzik, M. and Bork, P. (1997) *Cell* **88**, 9–11

Received 23 April 1997

## The Factor I and follistatin domain families: the return of a prodigal son

The Factor I/membrane-attack complex (FIMAC) domain has been identified in three complement proteins of immune defence and was believed to be unique [1]. It is present in the essential regulatory protease Factor I (FI) that is specific for the complement proteins C3b and C4b when complexed with their cofactors [2], and also in the terminal components C6 and C7 of the membrane-attack complex [3,4]. FI contains an N-terminal FIMAC domain (67 residues), a CD5 or scavenger-receptor cysteine-rich domain, two low-density-lipoprotein-receptor domains, and a serine-protease domain. Both C6 and C7 contain two C-terminal FIMAC domains. Here we show that database searches based on the FIMAC sequence indicated similarities with the follistatin (FS) sequence family (Figure 1). Similarities were also found in the predicted secondary structures of the FIMAC and FS domains (Figure 2). We therefore propose that the FIMAC domain is a distant member of the larger FS superfamily [5,6]. For FI, all the domain types in its structure have now been identified as members of larger cellular-receptor or extracellular-matrix domain superfamilies [7–9]. The FS superfamily of the extracellular matrix and the endocrine system now includes the immune equivalent of this domain.

Seven FIMAC sequences [3,4,10–12] were used in database searches of more than 59000 SWISSPROT sequences using BLITZ which is based on the program MPsrch with the BLOSUM62 amino acid substitution matrix [13]. Table 1 indicates hits with agrin in the FS family with 19–24% sequence matches. Lower matches in the top 50 hits (Table 1) included



**Figure 1 Occurrence of FIMAC and FS domains in the sequences of plasma, regulatory and extracellular matrix proteins**

Each FIMAC or FS domain is denoted by a box.

ACCESSION	PROTEIN	-----	-----	-----	-----	-----	-----	-----	-----
CFAL_HUMAN	FI (Human)	LSCDKVF.....	CQFWQRCIEGT.....	CV.C.KLPYQCPK...NGTAVCATNR...	RSFPTYCQKSLVCECLHPGT.....	KFLNNGICT			
XLC3BC4B	FI (Xenopus)	LSCHKVF.....	CAPWQRCVAGV.....	CR.C.KLPYQCPK...NATTEVCTDGG...	RKLQSYCQLKSVCECNPLNSK.....	YRFSSEAPCT			
MMU47810	FI (Mouse)	RSCNKVF.....	CQFWQRCIEGT.....	CI.C.KLPYQCPR...AGTPVCAMNG...	RSYPTYCHQKSFCECLHPE.....	IKFHSNAGTCA			
CO6_HUMAN	C6 rpt 1	LTCLKGH.....	CQLGQKQSGSE.....	CI.CMSPEEDCSH...HSEDLVCFDTSNDNYFTS	PACKFLAEKCLLNQO.....	LHFLHIGSCQ			
CO6_HUMAN	C6 rpt 2	ESCGYDT.....	CYDWEKCSASTSK.....	CV.C.LLPPQCFK...GGNLYCVKMGSSSTSEKTLNICEVGT	IRCANRK.....	MEILHFGKCL			
CO7_HUMAN	C7 rpt 1	LTQAVPK.....	CQRWEKLNQSR.....	CV.C.KMPYECGP...SLDVCAQDERSKRILPLTVCKMHLVHCQGRN...		YTLTGRDSCA			
CO7_HUMAN	C7 rpt 2	KACGA.....	CPLWKGKDAESSK.....	CV.C.REASECEE...EGFSICVEVNG...KEQTMSECEAGALRCRGS.....		SVTISRPT			
SPRC_HUMAN	SPARC	NPCQNH.....	CKHGKVCLEENNTPM.....	CV.C.QDPTSCPA.PIGEPEKVCNSDN...	KTFDSSCHFFATKCTLEGTKKGGHK.....	LHLDYIGPCK			
SPRC_CAEL	SPARC	NPCEDHQ.....	CGWGRECVVGGKGEPT.....	CE.C...ISKPELDDGDPMDKVCANNN...	QTFTSLCDLYRERCLCKRKSKECSKAFNAKVHLEYLGECK				
HSHEVIN	HEVIN (Human)	DSCMSFQ.....	CKRGHICKADQGGKPH.....	CV.C.QDPVTCPP...TKPLDQVCGTDN...	QTYASSCHLFFATKCRLEGTKKGGHQ.....	LQLDYFGACK			
AGRI_CHICK	AGRIN rpt 1	DACRGM.....	CGFGAVCERSPTDPSQAS.....	CV.C.KKTACPV...VVAPVCGSDY...	STYSNECELEKACQCNQQR...	IKVISKPGCC			
AGRI_CHICK	AGRIN rpt 2	DPCAENV.....	CSFGSTCVRSADGGTAG.....	CV.C...PASCSSG...VAESIVCGSDG...	KDYRSECCLNKHACDKQEN...	VFKKFDGACD			
AGRI_CHICK	AGRIN rpt 3	PCKGIL.....	NDMNRVCRVNPRTTRV.....	EL.L.SRPENCSS...KREPVCDDG...	VTYASECVMGRTGAIIRGLE...	IQKVRSGCQC			
AGRI_CHICK	AGRIN rpt 4	DKCKDE.....	CKFNAVCLKRWHAR.....	CS.C.DRITCDG...TYRPPVCARDS...	RTYSNDCERQNAKACHQKAA...	IPVKHSGPCD			
AGRI_CHICK	AGRIN rpt 5	SPCLSV.....	CTFGATCVVKNREP.....	CE.C...QVQVCG...RYDPVCGSDN...	RTYGNPCELNMAACVLRKRE...	IRVKHKGPCD			
AGRI_CHICK	AGRIN rpt 6	RCGK.....	CQFGAICEAETGR.....	CV.C.PTECV...SSQPVCVGTG...	NTYGSECELHVRACVQKN...	ILVAQGDCK			
AGRI_CHICK	AGRIN rpt 7	SCGTTV.....	CSFGSTCVGGQ.....	CV.C...PRCEQ...QPLAQVCGTDG...	LYYDNRELRAASCCQKKS...	IEVAKMGPE			
AGRI_CHICK	AGRIN rpt 8	DECGSGGSGSDGSECEQDR.CRHYGWDEDAEDDRVC.C.	DFTCLA...VPRSVCVCGSDD...	VTYANECLEKKTREKRON...		LYTVSQAGR			
AGRI_CHICK	AGRIN rpt 9	KSCCEMS.....	CEFGATCVVNGFAH.....	CE.C...PSPLCSE...ANMTKVCVCGSDG...	VTYGDCQLKTIACRQGL...	ITVKVHGACH			
FSA_HUMAN	FS rpt 1	ETCENVD.....	CGPGKCKRMNKNKPR.....	CV.C...APDCSN...ITWKGVPVCLDG...	KTYRNECALLKARCKEPE...	LEVQYQGRCK			
FSA_HUMAN	FS rpt 2	KTCRDVF.....	CPGSGTVCVQDTNNAY.....	CVTC...NRICPE.PASSEQYLCGNDG...	VTYSSACHLRKATCLLGRS...	IGLAYEGKCI			
FSA_HUMAN	FS rpt 3	KSCEDIQ.....	CTGGKCKLWDFKVGRR.....	CSLC...DELCPD...SKSDEPVCASDN...	ATYASECAMKKAACSSGVL...	LEVKHSGCCN			
S51362	FRP (Mouse)	KICANVF.....	CGAGRECAVTEKGEPT.....	CL.C...IEQCKP...HKRPVCGSNG...	KTYLNHCELHRDACLITGSK...	IQVDYDGHCK			
		5	10 15 20	25	30 35 40 45	50 55 60	65 70		
RESIDUE CONSERVATION (52)		-----	-----	-----	-----	-----	-----		
>90% conservation	C	C	C C	C C	C C	VCG	Y C C	L G C	
>70% conservation	C	C	C G C	CV C	C	VCG D	TY C L C	L L Y G C	
PREDICTED STRUCTURE (52)									
GOR I	ttttttt	tttttEttttttt	tE E	ttttttt	tttEEEEtttt	tttttHHHHHHHHHHHHHHH	EEEEtttttt		
GOR III	ctttttt	cttttEttttccct	EE c	ctttccc	ccctcEEEEctt	ccccttHHHHHHHHHHHt*H	EEEEtcccc		
Chou-Fasman	ttttEEE	ttttHEEttttEE	EE E	ttttttt	*tttEEEEtttt	*EEtHHHHHHHHHHHHHHH	E*H*tttttt		
SAPIENS	ooiooio	ioiooHHHooiooooo	ii i	ooiooio	ioooooioiooo	EEEOHHHHHHHHHHHHHHH	EEEEoioio		
PHD	lllEEll	lllllEEEElllllll	EE l	lllllll	lllllllEElll	llllllHHHHHHlllllll	EElllElll		
Averaged Structure		EEEE	EE E		EEEE	EE HHHHHHHHHHHHHH	EEEE		
PREDICTED ACCESSIBILITY (52)									
Hydropathy	ooiooio	ioiooioioooooii	ii i	oiooio	ioooooiiiooi	ooiooioiooioiooooo	ioiooioio		
PHD solvent access	ooiooi	i.ioooi.iooooooo	ii i	o.oioo	oo.ooioiooo	.ioioioiooioiooooo	iiioiooo.o		
SAPIENS solvent access	ooiooio	ioiooiooooooo	ii i	ooiooio	ioooooiiiooo	ooiooioiooioiooooo	ioooooioio		
Averaged Accessibility	ooiooio	ioiooioiooooooo	ii i	ooiooio	ioooooiiiooo	ooiooioiooioiooooo	ioiooioio		

**Figure 2 Summary of the multiple sequence alignment of FIMAC/FS sequences**

A total of 52 was used in the full alignment. The 29 not shown in the Figure have SWISSPROT accession codes SPRC\_BOVIN, SPRC\_MOUSE, SPRC\_RAT, SPRC\_CHICK, SPRC\_XENLA, QR1\_COTJA, SC1\_RAT, DYAGGR, AGRI\_RAT, FSA\_PIG, FSA\_RAT, FSA\_SHEEP, FSA\_XENLA, from which the full alignment is readily reconstructed. The two Cu<sup>2+</sup>-binding regions of the SPARC sequences are in **bold**. The alignment yielded a consensus length of 74 residues which is conserved in over 50% of the 52 sequences. The most commonly occurring residues that show 90% and 70% conservation within a given subtype (I=V=L=M; D=E; R=K=H; F=Y=W=H; G=A=S) are indicated. The averaged predicted structures are based on the presence of at least 2/5 secondary structures or 2/3 accessibility states at each residue position. Abbreviations are as follows; H,  $\alpha$ -helix; E,  $\beta$ -sheet; l, c and t, loop, coil and turn; o, solvent-accessible; i, solvent-inaccessible.

**Table 1 BLITZ database search of more than 59 000 sequences**

Probe FIMAC sequence (see Figure 2)	FS sequences detected ( <i>P</i> better than 10 <sup>-4</sup> )	Probability ( <i>P</i> )	Total of FS sequences in top 50 positions
CFAL_HUMAN (residues 41–107)	AGRI_RAT	4 × 10 <sup>-5</sup>	10
XLC3BC4B (residues 39–101)	AGRI_RAT	4 × 10 <sup>-6</sup>	9
MMU47810 (residues 44–110)	AGRI_CHICK	3 × 10 <sup>-7</sup>	16
	AGRI_RAT	9 × 10 <sup>-7</sup>	
CO6_HUMAN (residues 766–837)	None		17
CO7_HUMAN (residues 771–839)	AGRI_CHICK	7 × 10 <sup>-8</sup>	16
	AGRI_RAT	9 × 10 <sup>-6</sup>	

other agrin repeats, FS, SPARC (secreted protein, acidic, rich in cysteine; also known as osteonectin) and the SPARC homologues QR1 and SC1, all of which belong to the FS family. The BLITZ matches with the FS domains of chicken and rat agrin extended over 52–65 residues of the probe FIMAC sequences (Table 1) and aligned eight to ten of the eight or ten cysteine residues in the FIMAC sequences. Another database search using BLASTP [14] to scan more than 71 000 sequences showed that AGRI\_RAT and AGRI\_CHICK were scored with statistically significant probabilities (*P*) of 10<sup>-5</sup>–10<sup>-6</sup> by FIMAC sequences from FI and C7. Likewise version 3 of FASTA [15] yielded similar results with more than 59 000 sequences, in which FS domains in agrin, QR1 and FS also scored highly with probabilities *E*(59 000) less than 0.05, and lower matches included other FS domains. FASTA matched eight to ten of the eight or ten cysteine residues of FIMAC with their equivalents in the FS domain. Blockmaker

[16], with the FIMAC sequences and the FS domains from agrin and SPARC (Figure 2), produced two blocks, one of which aligned with both the FIMAC and FS sequences. Likewise, MACAW [17] generated a 28-residue block of ten aligned sequences from the top matches in the BLITZ search, which corresponded to that identified by Blockmaker [16] (residues 31–60 in Figure 2). Whereas the 28-residue block lacked the C6 or C7 FIMAC sequences, its occurrence was statistically significant.

Since these analyses showed a relationship between the FIMAC and FS sequences, a combined alignment was constructed using 52 sequences (Figure 2). The consensus length is 74 residues with ten conserved cysteine residues. These cysteine residues are assumed to be bridged. Although the disulphide pairings are unknown, cysteine mutations in the first FIMAC of C6 and C7 suggest that Cys<sup>3</sup>–Cys<sup>14</sup> are paired (the first and third Cys in

Figure 2). Cys<sup>8</sup>, Cys<sup>24</sup>, Cys<sup>26</sup> and Cys<sup>59</sup> (the second, fourth, fifth and ninth Cys in Figure 2) are missing in agrin domain 3 and may form two pairs if the domain structure is unaltered. Whereas C-terminal sequence similarities with the disulphide-rich Kazal-type inhibitors of the ovomucoid family and N-terminal similarities with the epidermal-growth-factor family have been noted [6,8], the proposed disulphide pairings in the alignment are mutually exclusive in the two families [5]. Figure 2 also shows conserved buried hydrophobic residues. These are attributable to the packing of the protein core, and most are present in both the FIMAC and FS sequences.

The FIMAC and FS sequences were also compared by computing consensus secondary-structure predictions from the alignment [18,19]. Use of the GORI, GORIII, Chou Fasman, PHD and SAPIENS prediction algorithms gave an averaged  $\beta\beta\beta\alpha\beta$  structure with five  $\beta$ -strands and one  $\alpha$ -helix (Figure 2). Residues 38–74 resemble the Kazal-type inhibitors. Interestingly, they resulted in a  $\beta\beta\alpha\beta$  prediction that is very similar to the observed  $\beta\beta\alpha\beta$  secondary structure in ovomucoid when analysed using DSSP (Brookhaven database codes: 1ovo-4ovo). Cys<sup>35</sup> and Cys<sup>38</sup> are located on the  $\alpha$ -helix of ovomucoid and correlate well with Cys<sup>52</sup> and Cys<sup>59</sup> in FIMAC/FS, which are located on the  $\alpha$ -helix predicted between residues 51 and 64 (Figure 2). Cys<sup>32</sup> and Cys<sup>59</sup> in FIMAC/FS will be positioned on the same side of this  $\alpha$ -helix as Cys<sup>35</sup> and Cys<sup>38</sup> in ovomucoid, since Cys<sup>52</sup> and Cys<sup>59</sup> will be separated by an extra turn of the predicted  $\alpha$ -helix in the FIMAC/FS domain compared with ovomucoid. The presence of the extra N-terminal residues in FIMAC before the region of similarity with ovomucoid implies that a structural relationship to ovomucoid is possible, but this will be modified. Figure 2 also shows large sequence insertions in the alignment, and these occur in regions predicted to be surface loops as desired. The loop between residues 64 and 65 corresponds to the Cu<sup>2+</sup>-binding region of the SPARC proteins [5] that is implicated in cellular proliferation, but is absent from other FIMAC and FS domains. More importantly, application of averaged secondary-structure predictions to each of the FIMAC and FS sequence families yielded results that were very similar to the predicted  $\beta\beta\beta\alpha\beta$  structure for all 52 sequences and support the proposed identity between the two families.

The phylogenic relationship between the FIMAC and FS sequences was investigated using PHYLIP [20]. An unrooted tree showed that the FIMAC sequences occupied a separate branch from the FS sequences. The lengths of the exons in SPARC, agrin, FI, C6 and C7 are in agreement with Figure 2, whereas their boundaries are not conserved. In murine SPARC, the FS domain is encoded by exons 5 and 6 with intron boundaries of class 1-1 [5], and correlates well with the FS domains in agrin, which are encoded by one or two exons and have intron boundaries of class 1-1 [21]. In human FI, the FIMAC domain is within exon 2, but with intron boundaries of class 0-1 [22]. The two FIMAC domains in human C6 and C7 are encoded across exons 15, 16 and 17 with intron boundaries of class 1-2, 2-1 and 1-undefined respectively [3,4]. The lack of conserved intron/exon structure is similar to that found in the serine-protease domain [22].

The relationship between the FIMAC and FS families is matched functionally in that all proteins containing FIMAC and FS domains are extracellular and participate in protein-protein interactions (Figure 1). The FS domain has been identified in extracellular matrix proteins that modulate cell-matrix interactions (SPARC), induce aggregation of nicotinic acetylcholine receptors (agrin) and in ovaries and the pituitary that bind cytokines (FS) [5,6]. Hevin is isolated from high endothelial venules of tonsils which allow high levels of lymphocyte extra-

vasion from blood and may facilitate this lymphocyte migration [23]. SPARC is released by platelet degranulation, is synthesized by fibroblasts and macrophages at sites of wound repair, and may regulate deposition or assembly of extracellular matrix proteins. The specificity of the FS domain for its ligand is indicated by SPARC and FS. SPARC binds to platelet-derived growth factor, albumin, thrombospondin and various collagen types. FS binds to activin and inhibin, which are transforming-growth-factor- $\beta$ -like cytokines. Despite the sequence similarity to ovomucoid, no protease-inhibitory activity has been reported to date for FS-containing proteins [6]. FI interacts with the complement components C3b and C4b, whereas C6 and C7 interact with complement component C5b during formation of the membrane-attack complex. As C3b, C4b and C5b are all related in sequence, it will be of interest to determine whether all three contain a similar target fold for FIMAC. Given the relationship to ovomucoid, it will be of interest to determine whether FIMAC in FI can inhibit its own serine-protease domain. From X-ray and neutron-scattering analyses of FI, one model that is consistent with the data is a semi-compact V-shaped structure, the dimensions of which place these two domains proximate to each other [24].

We thank the Wellcome Trust for grant support and Dr. R. B. Sim for useful discussions.

Christopher G. ULLMAN and Stephen J. PERKINS<sup>1</sup>

Department of Biochemistry and Molecular Biology, Royal Free Hospital School of Medicine, Rowland Hill Street, London NW3 2PF, U.K.

<sup>1</sup> To whom correspondence and requests for offprints should be addressed.

- Law, S. K. A. and Reid, K. B. M. (1995) Complement, 2nd edn., IRL Press, Oxford
- Sim, R. B., Day, A. J., Moffatt, B. E. and Fontaine, M. (1993) *Methods Enzymol.* **223**, 13–35
- Hobart, M. J., Fernie, B. and DiScipio, R. G. (1993) *Biochemistry* **32**, 6198–6205
- Hobart, M. J., Fernie, B. A. and DiScipio, R. G. (1995) *J. Immunol.* **154**, 5188–5194
- Lane, T. F. and Sage, E. H. (1994) *FASEB J.* **8**, 163–173
- Patthy, L. and Nikolics, K. (1993) *Trends Neurosci.* **16**, 76–81
- Perkins, S. J. and Smith, K. F. (1993) *Biochem. J.* **295**, 109–114
- Moestrup, S. K. (1994) *Biochim. Biophys. Acta* **1197**, 197–213
- Resnick, D., Pearson, A. and Krieger, M. (1994) *Trends Biochem. Sci.* **19**, 5–8
- Catterall, C. F., Lyons, A., Sim, R. B., Day, A. J. and Harris, T. J. (1987) *Biochem. J.* **242**, 849–856
- Kunnath-Muglia, L. M., Chang, G. H., Sim, R. B., Day, A. J. and Ezekowitz, R. A. (1993) *Mol. Immunol.* **30**, 1249–1256
- Minta, J. O., Wong, M. J., Kozak, C. A., Kunnath-Muglia, L. M. and Goldberger, G. (1996) *Mol. Immunol.* **33**, 101–112
- Sturrock, S. S. and Collins, J. F. (1993) MPsrch, Version 1.5, Biocomputing Research Unit, University of Edinburgh
- Altschul, S. F., Warren, G., Miller, W., Myers, E. W. and Lipman, D. J. (1990) *J. Mol. Biol.* **215**, 403–410
- Pearson, W. R. and Lipman, D. J. (1988) *Proc. Natl. Acad. Sci. U.S.A.* **85**, 2444–2448
- Henikoff, S., Henikoff, J. G., Alford, W. J. and Pietrokovski, S. (1995) *Gene* **163**, GC (Gene-COMBIS) 17–26
- Lawrence, C. E., Altschul, S. F., Boguski, M. S., Liu, J. S., Neuwald, A. F. and Woolton, J. C. (1993) *Science* **262**, 208–214
- Edwards, Y. J. K. and Perkins, S. J. (1996) *J. Mol. Biol.* **260**, 277–285
- Brissett, N. C. and Perkins, S. J. (1996) *FEBS Lett.* **388**, 211–216
- Felsenstein, J. (1989) *Cladistics* **5**, 164–166
- Rupp, F., Özçelik, T., Linial, M., Peterson, K., Francke, U. and Scheller, R. (1992) *J. Neurosci.* **12**, 3535–3544
- Vyse, T. J., Bates, G. P., Walport, M. J. and Morley, B. J. (1994) *Genomics* **24**, 90–98
- Girard, J. P. and Springer, T. A. (1995) *Immunity* **2**, 113–123
- Perkins, S. J., Smith, K. F. and Sim, R. B. (1993) *Biochem. J.* **295**, 101–108