

Polysaccharide Hydrolase Folds Diversity of Structure and Convergence of Function

MICHAEL E. HIMMEL,^{*,1} P. ANDREW KARPLUS,²
JOSHUA SAKON,² WILLAM S. ADNEY,¹ JOHN O. BAKER,¹
AND STEVEN R. THOMAS¹

¹*Biotechnology Center for Renewable Fuels and Chemicals,
National Renewable Energy Laboratory, Golden, CO 80401;*
and ²*Section of Biochemistry, Molecular and Cell Biology
Cornell University, Ithaca, NY 14853*

ABSTRACT

Polysaccharide glycosyl hydrolases are a group of enzymes that hydrolyze the glycosidic bond between carbohydrates or between a carbohydrate and a noncarbohydrate moiety. Here we illustrate that traditional schemes for grouping enzymes, such as by substrate specificity or by organism of origin, are not appropriate when thinking of structure–function relationships and protein engineering. Instead, sequence comparisons and structural studies reveal that enzymes with diverse specificities and from diverse organisms can be placed into groups among which mechanisms are largely conserved and insights are likely to be transferrable. In particular, we illustrate how enzymes have been grouped using protein sequence alignment algorithms and hydrophobic cluster analysis. Unfortunately for those who seek to improve cellulase function by design, cellulases are distributed throughout glycosyl hydrolase Families 1,5,6,7,9, and 45. These cellulase families include members from widely different fold types, i.e., the TIM-barrel, $\beta\alpha\beta$ -barrel variant (a TIM-barrel-like structure that is imperfectly superimposable on the TIM-barrel template), β -sandwich, and α -helix circular array. This diversity in cellulase fold structure must be taken into account when considering the transfer and application of design strategies between various cellulases.

*Author to whom all correspondence and reprint requests should be addressed.

Index Entries: Cellulases; xylanases; amylases; glycosyl hydrolases; structural folds; X-ray structures; hydrophobic cluster families.

INTRODUCTION

In nature, the enzymatic degradation of cellulose is a fundamental mechanism of biomass conversion and carbon cycling in the biosphere. To produce alcohol fuels from lignocellulosic biomass, however, this process must be understood at the molecular level to develop highly efficient and cost-effective catalysts. One important step toward this goal is to determine key structure–function relationships of enzymes displaying activity on water-insoluble substrates, such as cellulose and other plant polysaccharides. In order to define and understand complex molecular mechanisms, detailed structural information, such as that determined from NMR or X-ray diffraction studies, is essential.

Understanding the evolutionary and mechanistic relationships of enzymes that catalyze similar reactions is, therefore, a highly desirable objective for those who design new strategies to improve enzyme function by site-directed mutagenesis (SDM). An examination of the SWISS-PROT data base by Orengo and coworkers (1) using sequence alignments (FASTA) revealed that the ~28,000 entries were reducible to no less than about 7700 unambiguously related groups. Groupings with such limited sensitivity tell us little about enzyme relatedness and are not a serviceable tool for understanding enzyme design. A more sensitive method is hydrophobic cluster analysis (HCA), which relies on the basic rules that underlie the folding of globular proteins and uses a two-dimensional plot to display the amino acid sequence of a protein depicted as an “unrolled longitudinal cut” of a cylinder (2). The “helical net” produced by this graphical display allows the full sequence environment of each amino acid to be examined. HCA has been an extremely powerful method for classifying glycosyl hydrolases (2). Gilkes and coworkers (3) originally proposed nine families, based on the glycosyl hydrolase sequences available at that time, and in the ensuing five years, these researchers have added substantially to the original classification list (4,5). Development of glycosyl hydrolase classification is shown chronologically in Fig. 1. Today, Bairoch (6) has identified 56 families. This classification system provides a powerful tool for glycosyl hydrolase enzyme engineering studies, because many enzymes critical for industrial processes have not yet been crystallized or subjected to structure analysis. This article will review the correlations between polysaccharidase function and fold structure based on existing family assignments and reported macromolecular structures.

Enzyme Structure

Protein domains are grouped into four general structural categories (all- α , all- β , $\alpha + \beta$, and α/β) (7,8). Proteins of the all- α class are usually comprised of multiple α -helices that may be oriented either along a common

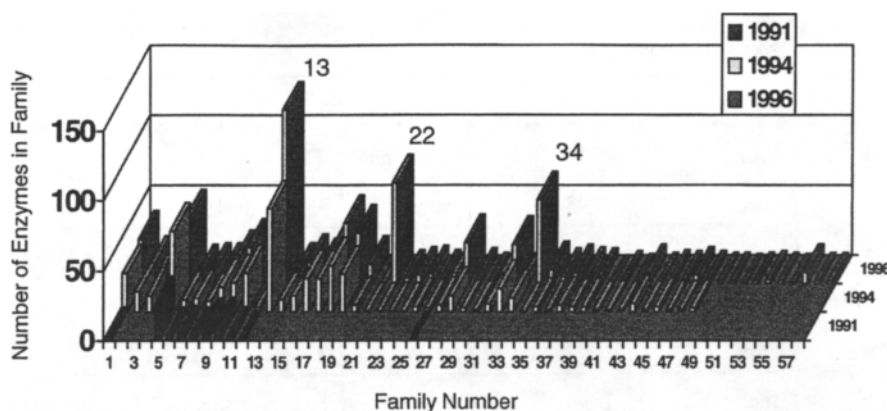


Fig. 1. The chronological development of O-glycosyl hydrolase classification by protein sequence alignment and hydrophobic cluster analysis (3,4,6).

bundle axis or randomly (9). Proteins of the all- β class contain β -strands that can be oriented either parallel or antiparallel, or in a mixture of the two (10). The $\alpha + \beta$ and α/β categories are distinguished by considering that α/β proteins have alternating β -strand and α -helical segments, whereas $\alpha + \beta$ proteins tend to contain regions definable as "mostly α " and "mostly β " (11). A common example of the α/β class is the TIM-barrel, named after the archetype of this fold, triose phosphate isomerase. In TIM-barrel proteins, the internal barrel is comprised of eight parallel β -strands, and the outer shell contains eight α -helices oriented with a cant relative to the axis of the barrel. Some protein domains do not fall into any of these categories and are grouped as irregular folds. Proteins representative of these domain (or fold) classes are myoglobin (all- α helix), immunoglobulin (all- β strand), cytochrome b5 ($\alpha + \beta$), and triose phosphate isomerase (α/β) (7).

It is inferred that all proteins that have recognizable sequence similarity will have the same fold type. In many cases, the fold will be unique to that single family of proteins; such folds are known as structural singlets (1). In other cases, a domain structure (fold) may be shared by two or more proteins that appear unrelated by sequence and function. Such folds have been termed superfolds. Orengo and coworkers identified nine protein superfolds that were the α/β -doubly wound fold, the TIM-barrel, the split α/β -sandwich, the greek-key immunoglobulin fold, the up-down four-helix bundle, the globin fold, the β -jelly roll, the β -trefoil, and the ubiquitin $\alpha\beta$ -roll (1).

Structure of Polysaccharide Hydrolases

A recent survey of the 3000 or more protein structures in the Brookhaven National Laboratory Protein Data Bank (PDB) revealed that only 6 cellulases, 8 xylanases, and 10 α -amylases, β -amylases, and glucoamylases have been deposited. A seventh cellulase structure was recently

Table 1
Distribution of Cellulases Among Glycosyl Hydrolase Families

Family	PDB file	Year filed	Res., Å	Fold class	Fold*	Ref.
Family 1	1cbg	1994	2.2	beta/alpha	TIM-barrel	16
Family 5	1cec	1995	2.2	beta/alpha	TIM-barrel	17
	1ece	1996	2.4	beta/alpha	TIM-barrel	12
Family 6	3cbh	1990	2	beta/alpha	TIM-barrel variant	18
	1tml	1993	1.8	beta/alpha	TIM-barrel variant	19
Family 7	1cel	1994	1.8	all beta	beta-sandwich, 12-14 strands	20
Family 9	1clc	1995	1.9	all alpha	6 alpha-helices/circular array	21
Family 45	1eng	1993	1.6	all beta	closed barrel, Barwin-like	22

*Assignment made by SCOP (15).

resolved for the endoglucanase EI from *Acidothermus cellulolyticus* (12). For the most part, the structural information available pertains to the catalytic domain only; however, two structures for cellulase cellulose binding domains (CBDs) have been filed with PDB (1cbh and 1exg). These structures are defined as small, all- β strand domains with 7 or 8 strands per molecule and were solved using 2D NMR techniques (13,14). This list will, of course, increase with time. Information regarding the structures of polysaccharide glycosyl hydrolases is shown in Tables 1–3. The structural classification given in these tables is based on recommendations from Structural Classification of Proteins (SCOP) (<http://www.pdb.bnl.gov/scop>) (15) and the recent glycosyl hydrolase family information taken from Bairoch, <http://www.expasy.hcuge.ch/cgi-bin/lists?glycosid.txt> (6).

Tertiary structure and key residues at active sites are generally better conserved than amino acid sequence, so it is no surprise that structural studies, combined with sequence comparisons directed at active site residues, have allowed many families to be grouped in clans that have a common fold and a common catalytic apparatus (39). Five such clans recently proposed for the glycosyl hydrolases are GH-A (including Families 1, 2, 5, 10, 17, 30, 35, 39, and 42); GH-B (Families 7 and 16); GH-C (Families 11 and 12); GH-D (Families 27 and 36); and GH-E (Families 33 and 34) (6). Among these, Clans GH-A and GH-B include cellulases.

DISCUSSION

The objective of this article is to compare the gross structural features recently made available for selected polysaccharide glycosyl hydrolases and to draw correlations with function. To accomplish this, the

Table 2
Distribution of Xylanases Among Glycosyl Hydrolase Families

Family	PDB file	Year file	Res., Å	Fold class	Fold*	Ref.
Family 10	1xas	1994	2.6	beta/alpha	TIM-barrel	23
	1xyz	1995	1.4	beta/alpha	TIM-barrel	24
	1xys	1994	2.5	beta/alpha	TIM-barrel	25
	2exo	1994	1.8	beta/alpha	TIM-barrel	26
Family 11	1xyp	1994	1.5	all beta	beta-sandwich, 12-14 strand	27
	1xyn	1994	2	all beta	beta-sandwich, 12-14 strand	27
	1xnd	1994	1.8	all beta	beta-sandwich, 12-14 strand	28
	1bcx	1994	1.8	all beta	beta-sandwich, 12-14 strand	29

*Assignment made by SCOP (15).

Table 3
Distribution of Starch-Degrading Enzymes Among Glycosyl Hydrolase Families

Family	PDB file	Year filed	Res., Å	Fold class	Fold*	Ref.
Family 13	2taa	1982	3	beta/alpha	TIM-barrel	30
	1cgt	1993	2	beta/alpha	TIM-barrel	31
	1cyg	1993	2.5	beta/alpha	TIM-barrel	32
	1ppi	1994	2.2	beta/alpha	TIM-barrel	33
	2aaa	1991	2.1	beta/alpha	TIM-barrel	34
	1amg	1993	2.2	beta/alpha	TIM-barrel	35
	1amy	1994	2.8	beta/alpha	TIM-barrel	36
Family 14	1byb	1994	1.9	beta/alpha	TIM-barrel	37
Family 15	3gly	1994	2.2	beta/alpha	6 alpha-helices/circular array	38

*Assignment made by SCOP (15).

families of glycosyl hydrolases classified by HCA and sequence alignment have been examined. Figure 2 shows the general substrate specificities of enzymes assigned to these families. Several families are highly conserved relative to substrate preference; these include the β -galactosidases (Family 2), β -glucosidases (Family 3), xylanases (Families 10 and 11), α -amylases (Family 13), β -amylases (Family 14), glucoamylases (Family 15), lichenases (Families 16 and 17), chitinases (Families 18 and 19), lysozymes (Families 21–24), neuraminidases (Families 33 and 34),

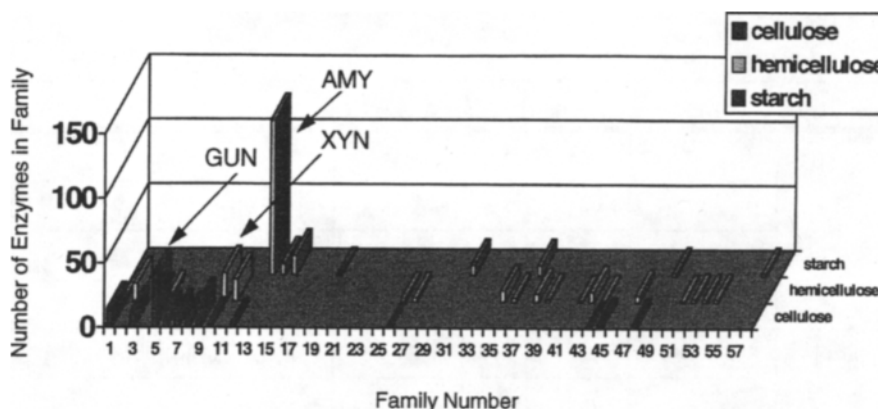


Fig. 2. The polysaccharide substrate specificity and distribution for families of O-glycosyl hydrolases.

and trehalases (Family 37). However, other families, such as Families 1, 4, 5, 39, 44, and others, include members that degrade different kinds of glycoside substrates.

Thus, within a single family, there may be diversity in the types of polysaccharides attacked, even though each family is expected to show a common tertiary structure and have a common mechanism of hydrolysis (4). Conversely, Tables 1–3 show that each general polysaccharide type may be attacked by hydrolases from various families, indeed with representation from widely different fold types. For instance, structurally known enzymes that degrade cellulose include two families with the TIM-barrel superfold (i.e., cyanogenic β -D-glucosidase [1cbg] from Family 1 and endoglucanase C from *Clostridium thermocellum* [1cec] and endoglucanase EI from *A. cellulolyticus* [1ece] from Family 5), one family with a modified TIM-barrel fold (i.e., cellobiohydrolase II from *Trichoderma reesei* [3cbh] and endoglucanase 2 from *Thermomonospora fusca* [1tml] from Family 6), one family with a β -strand sandwich fold (i.e., cellobiohydrolase I from *T. reesei* [1cel] from Family 7), one family with a six α -helix circular array fold (i.e., endoglucanase D from *C. thermocellum* [1c1c] from Family 9), and one family that has a Barwin-like, all β -strand closed barrel fold (i.e., endoglucanase V from *Humicola insolens* [1eng] from Family 45).

Table 2 shows a similar case for the eight xylanase structures reported, where one family with the TIM-barrel superfold (i.e., Family 10, xylanase A from *Streptomyces lividans* [1xas], xylanase Z from *C. thermocellum* [1xyz], xylanase A from *Pseudomonas fluorescens* var. *cellulosa* [1xys], and cellulase/xylanase from *Cellulomonas fimi* [2exo]) and a second family with a β -strand sandwich fold (i.e., Family 11, xylanase II from *T. reesei* [1xyp], xylanase I from *T. reesei* [1xyn], xylanase II from *Trichoderma harzianum* [1xnd], and xylanase II from *Bacillus circulans* [1bcx]) were

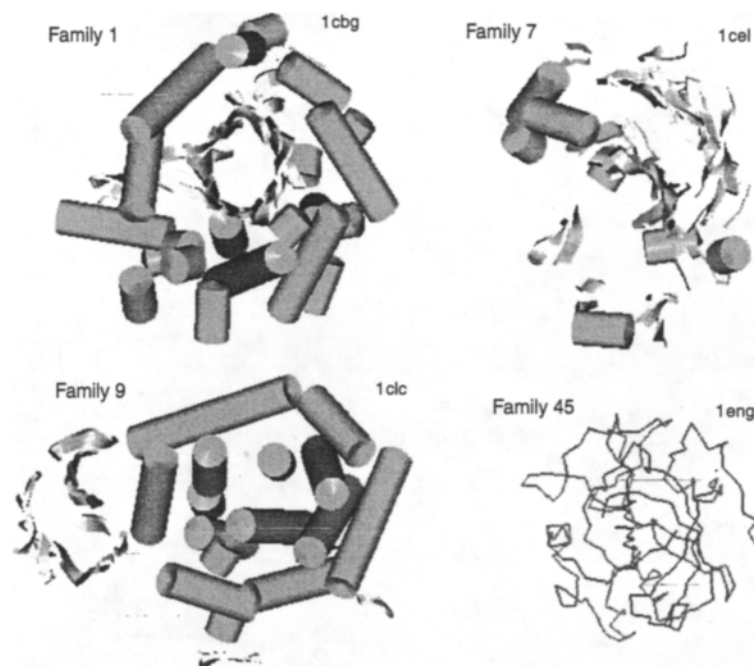


Fig. 3. Views of the secondary structures of cellulases and related enzymes from glycosyl hydrolase Families, 1, 7, 9, and 45. 1cbg is the cyanogenic-D-glucosidase from *Trifolia repens*, 1cel is the cellobiohydrolase I from *T. reesei*, 1clc is the endoglucanase D from *C. thermocellum*, and 1eng is the endoglucanase V from *H. insolens*. Secondary assignments were made by Kabish-Sander algorithms. The structures were generated from PDB files (Brookhaven National Laboratory) using Biosym Version 94, Biosym Technologies, San Diego, CA.

found. Views of secondary structures for cellulase and xylanase enzymes generated from PDB files and assigned to glycosyl hydrolase families are shown in Figs. 3–5.

Table 3 shows less structural diversity for the starch-degrading enzymes studied so far, because all members of both Families 13 and 14, the α -amylases and β -amylases, respectively, contain TIM-barrel super-folds. The Family 15 glucoamylase, an enzyme from *Aspergillus awamori*, has a fold type much different from the TIM-barrel, i.e., a six α -helix circular array also found for endoglucanase D from *C. thermocellum*. This suggests that Family 15 enzymes may be mechanistically more closely related to the Family 9 cellulases than to the other amylases. Whereas the TIM-barrel superfold is a highly versatile and robust structure that dominates in the collection of glycosyl hydrolases as a whole (1,5), clearly many other folds are represented.

The high levels of structural and mechanistic similarity that often occur between enzymes showing preferences for different substrates also

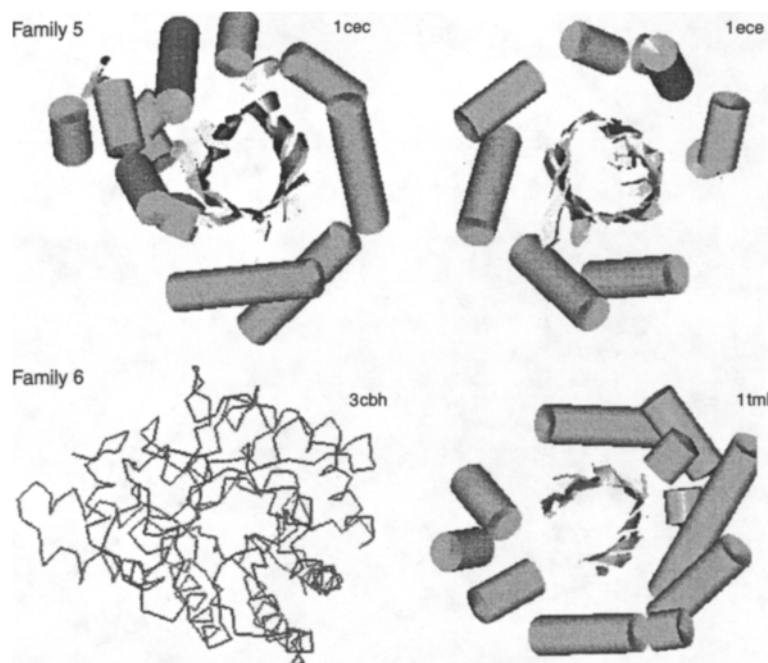


Fig. 4. Views of the secondary structures of cellulases from glycosyl hydrolase Families 5 and 6. 1cec is the endoglucanase C from *C. thermocellum*, 1ece is the endoglucanase EI from *A. cellulolyticus*, 3cbh is the cellobiohydrolase II from *T. reesei*, and 1tml is the endoglucanase 2 from *T. fusca*. The structures were generated from PDB files (Brookhaven National Laboratory) using Biosym Version 94, Biosym Technologies, San Diego, CA.

occur between enzymes from different environments and from different types of organisms. Within Family 5, for example, structural correlations were expected and have been recently confirmed experimentally between the endoglucanase C from *C. thermocellum* (1cec) and the endoglucanase EI from *A. cellulolyticus* (1ece). The α -carbon traces for these two enzymes are nearly superimposable (12). This is intriguing, considering that endoglucanase C was produced as part of a mesophilic cellulosomal cellulase system and EI from *A. cellulolyticus* is a highly thermal tolerant endoglucanase (optimum temperature of 81°C) derived from a hot spring bacterium (40). Furthermore, as shown by their grouping in the GH-A clan (6), the Family 1 enzymes and the Family 10 xylanases are structurally and mechanistic cousins to the Family 5 cellulases. Structural and functional similarity superceding taxonomic boundaries is also common. For example, the Family 11 xylanases are all β -strand sandwich folds (Table 2), yet three are fungal (*Trichoderma*) and one is from a bacterium (*Bacillus*). Even cursory examination of Fig. 5 confirms the high degree of structural similarity for the three xylanases shown from Family 11.

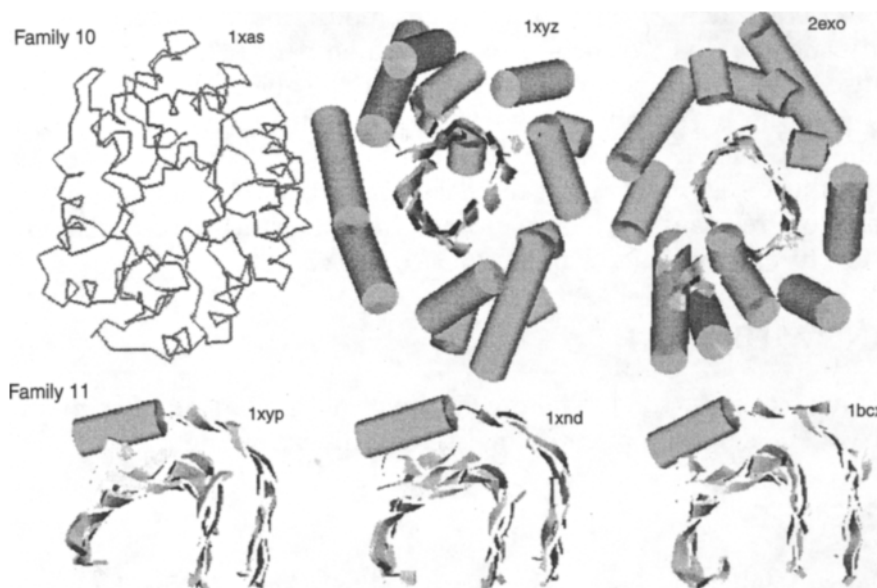


Fig. 5. Views of the secondary structures of xylanases from glycosyl hydrolase Families 10 and 11. 1xas is the xylanase A from *S. lividans*, 1xyz is the xylanase Z from *C. thermocellum*, 2exo is the cellulase/xylanase from *C. fimi*, 1xyp is the xylanase II from *T. reesei*, 1xnd is the xylanase II from *T. harzianum*, and 1bcx is the xylanase II from *B. circulans*. For clarity, structures for 1xys and 1xyn are not shown. The structures were generated from PDB files (Brookhaven National Laboratory) using Biosym Version 94, Biosym Technologies, San Diego, CA.

CONCLUSION

Ideally an experimental structure would be available for each glycosyl hydrolase, but we are limited by the size of the available structure data base at this time. Specifically, of the 811 glycosyl hydrolases listed by Bairoch (6), only 30 of those important for lignocellulosic biomass conversion have been subjected to X-ray diffraction and structure analysis. Still, nature has provided many independent structural schemes to bring a set of key active site amino acid residues into precise position to effect very similar or identical chemical reactions, in this case, transfer of the O-glycosyl bond to water. Thus, diversity of structure is a fact, and convergence of catalytic function is indicated for different types of glycosyl hydrolase folds. Glycosyl hydrolase Family 5 enzymes are especially important to researchers in the biomass conversion field, because one especially active and thermal tolerant endoglucanase, EI from *A. cellulolyticus*, belongs to Subclass 1 of this family. Although X-ray structures are currently available for only 2 of the 47 enzymes identified in Family 5 to date, endoglucanase C from *C. thermocellum* and endoglucanase EI from *A. cellulolyticus*, SDM work to improve Family 5 cellulases can proceed with a reasonable degree

of confidence that lessons learned from modifying the structure of one family member will translate to other members (41,42). In addition, because of the clear placement of Family 5 in the glycosyl hydrolase GH-A clan, insights based on studies of enzymes from glycosyl hydrolase Families 1, 2, 10, 17, 30, 35, 39, and 42 will also likely be relevant, even though those enzymes are not cellulases. In fact, such insights will probably be more relevant than those gleaned from the study of cellulases from the structurally unrelated Families 6, 7, 9, and 45.

ACKNOWLEDGMENT

This work was funded by the Biochemical Conversion Element of the Office of Fuels development of the US Department of Energy.

REFERENCES

- Orengo, C. A., Jones, D. T., and Thornton, J. M. (1994), *Nature* **372**, 631–634.
- Henrissat, B. (1991), *Biochem. J.* **280**, 309–316.
- Gilkes, N. R., Henrissat, B., Kilburn, D. G., Miller, R. C., Jr., and Warren, R. A. J. (1991), *Microbiol. Rev.* **55**, 303–312.
- Henrissat, B. and Bairoch, A. (1993), *Biochem. J.* **293**, 781–788.
- Henrissat, B., Callebaut, I., Fabrega, S., Lehn, P., Mornon, J.-P., and Davies, G. (1995), *Proc. Natl. Acad. Sci.* **92**, 7090–7094.
- Bairoch, A. (1996), SWISS-PROT Protein Sequence Data Bank (<http://expasy.hcuge.ch/cgi-bin/lists?glycosid.text>).
- Levitt, M. and Chothia, C. (1976), *Nature* **261**, 552–557.
- Efimov, A. V. (1994), *Structure* **2**, 999–1002.
- Harris, N. L., Presnell, S. R., and Cohen, F. E. (1994), *J. Mol. Biol.* **236**, 1356–1368.
- Chothia, C. and Janin, J. (1981), *Proc. Natl. Acad. Sci. USA* **78**, 4146–4150.
- Orengo, C. A. and Thornton, J. M. (1993), *Structure* **1**, 105–120.
- Sakon, J., Adney, W. S., Himmel, M. E., Thomas, S. R., and Karplus, P. A. (1996), *Biochemistry* **35**, 10648–10660.
- Kraulis, P. J., Clore, G. M., Nilges, M., Jones, T. A., Pettersson, G., Knowles, J., and Gronenborn, A. M. (1989), *Biochemistry* **28**, 7241.
- Xu, G.-Y., Ong, E., Gilkes, N. R., Kilburn, D. G., Muhandiram, D. R., Harris-Brandts, M., Carver, J. P., Kay, L. E., and Harvey, T. S. (1995), PDB entry 1exg.
- Murzin, A. G., Brenner, S. E., Hubbard, T., and Chothia, C. (1995), *J. Mol. Biol.* **247**, 536–540.
- Tolley, S. P., Barrett, T. E., Suresh, C. G., and Huges, M. A. (1993), *J. Mol. Biol.* **229**, 791.
- Dominguez, R., Souchon, H., Spinelli, S., Dauter, Z., Wilson, K. S., Chauvaux, S., Beguin, P., and Alzari, P. M. (1995), *Nat. Struct. Biol.* **2**, 569.
- Rouvainen, T., Rouvainen, J., Lehtovaara, P., Caldentey, X., Tomme, P., Claeysens, M., Pettersson, G., and Teeri, T. (1989), *J. Mol. Biol.* **209**, 167.
- Spezio, M., Wilson, D. B., and Karplus, P. A. (1993), *Biochemistry* **32**, 9906.
- Divne, C., Stahlberg, J., Reinikainen, T., Ruohonen, L., Pettersson, G., Knowles, J. K. C., Teeri, T. T., and Jones, T. A. (1994), *Science* **265**, 524.
- Alzari, P. M., Juy, M., and Souchon, H. (1993), *Biotechnol. Industrial Fermentation* **8**, 73.
- Davies, G. J., Dodson, G. G., Hubbard, R. E., Tolley, S. P., Dauter, Z., Wilson, K. S., Hjort, C., Mikkelsen, J. M., Rasmussen, G., and Schulein, M. (1993), *Nature* **365**, 362.
- Derewenda, U., Swenson, R., Green, R., Wei, Y., Morosoli, R., Shareck, F., Kluepfel, D., and Derewenda, Z. S. (1994), *J. Biol. Chem.* **269**, 20,811.

24. Dominguez, R., Souchon, H., Spinelli, S., Dauter, Z., Wilson, K. S., Chauvaux, S., Beguin, P., and Alzari, P. M. (1995), *Nat. Struct. Biol.* **2**, 569.
25. Harris, G. W., Jenkins, J. A., Connerton, I., Cummings, N., Lo Leggio, L., Scott, M., Hazlewood, G. P., Laurie, J. I., Gilbert, H. J., and Pickersgill, R. W. (1994), *Structure (Lond.)* **2**, 1107.
26. White, A., Withers, S. G., Gilkes, N. R., and Rose, D. R. (1994), *Biochemistry* **33**, 12,546.
27. Torronen, A. and Rouvinen, J. (1995), *Biochemistry* **34**, 847.
28. Campbell, R. L., Rose, D. R., Wakarchuk, W. W., To, R. J., Sung, W., and Yaguchi, M. (1994), PDB entry 1xnd.
29. Wakarchuk, W. W., Campbell, R. L., Sung, W. L., Davoodi, J., and Yaguchi, M. (1994), *Protein Sci.* **3**, 467.
30. Matsuura, Y., Kusunoki, M., Harada, W., and Kakudo, M. (1984), *J. Biochem. (Tokyo)* **95**, 697.
31. Klein, C. and Schulz, G. E. (1991), *J. Mol. Biol.* **217**, 737.
32. Kubota, M., Matsuura, Y., Sakai, S., and Katsube, Y. (1995), PDB entry 1cyg.
33. Qian, M., Haser, R., Buisson, G., Duee, E., and Payan, F. (1993), *J. Mol. Biol.* **231**, 785.
34. Boel, E., Brady, L., Brzozowski, A. M., Derewenda, Z., Dodson, G. G., Jensen, V. J., Petersen, S. B., Swift, H., Thim, L., and Woldike, H. F. (1990), *Biochemistry* **29**, 6244.
35. Morishita, Y., Matsuura, Y., Kubota, M., Sato, M., Sakai, S., and Katsube, Y. (1995), PDB entry 1amg.
36. Kadziola, A., Abe, J.-I., Svensson, B., and Haser, R. (1994), *J. Mol. Biol.* **239**, 104.
37. Mikami, B., Degano, M., Hehre, E. J., and Sacchettini, J. C. (1994), *Biochemistry* **33**, 7779.
38. Aleshin, A. E., Hoffman, C., Firsov, L. M., and Honzatko, R. B. (1994), *J. Mol. Biol.* **238**, 575.
39. Davies, G. and Henrissat, B. (1995), *Structure* **3**, 853–859.
40. Himmel, M. E., Adney, W. S., Grohmann, K., and Tucker, M. P. (1994), US Patent No. 5,275,944.
41. Wang, Q., Tull, D., Meinke, A., Gilkes, N. R., Warren, R. A. J., Aebersold, R., and Withers, S. G. (1993), *J. Biol. Chem.* **268**, 14,096–14,102.
42. Bortoli-German, I., Haiech, J., Chippaux, M., and Barras, F. (1995), *J. Mol. Biol.* **246**, 82–94.