

# Supplemental Data

## Extent and Origins of Functional Diversity in a Subfamily of Glycoside Hydrolases

### Authors:

Evan M. Glasgow<sup>1,2</sup>, Kirk A. Vander Meulen<sup>1,2</sup>, Taichi E. Takasuka<sup>1,2,3</sup>, Christopher M. Bianchetti<sup>1,2,4</sup>, Lai F. Bergeman<sup>1,2</sup>, Samuel Deutsch<sup>5</sup>, and Brian G. Fox<sup>1,2</sup>

### Affiliations:

<sup>1</sup>Great Lakes Bioenergy Research Center, Madison, WI, USA

<sup>2</sup>Department of Biochemistry, University of Wisconsin – Madison, USA

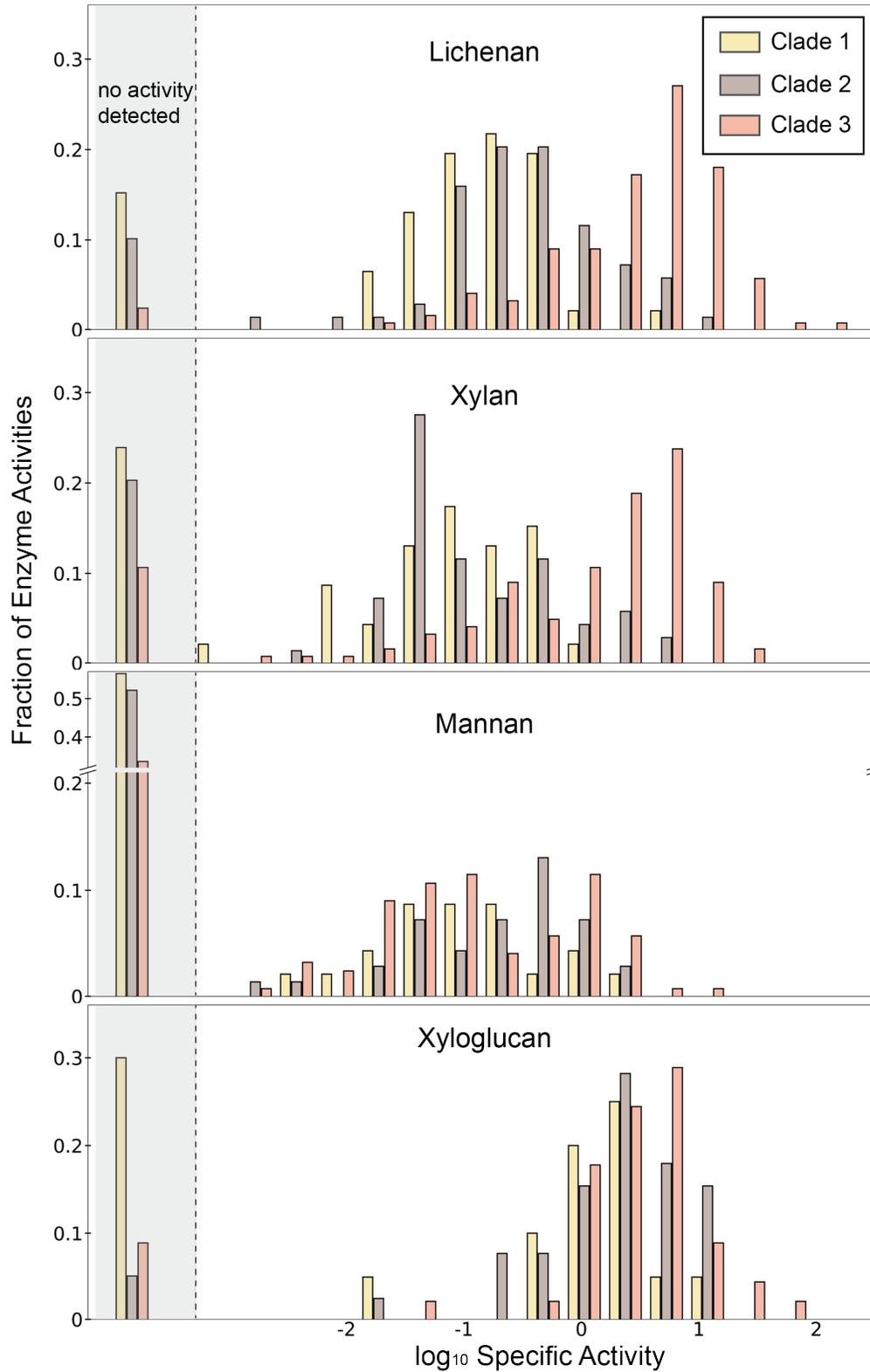
<sup>3</sup>Research Faculty of Agriculture, Hokkaido University, Sapporo, Japan

<sup>4</sup>Department of Chemistry, University of Wisconsin – Oshkosh, USA

<sup>5</sup>DOE Joint Genome Institute, Walnut Creek, CA, USA

*A Distinct Activity Distribution for Clade 3 Enzymes* — Supplementary Figure S1 illustrates that activity values in Clade 3 enzymes are shifted toward higher activities for lichenan and xylan. For mannanase activities, the assays are likely sampling from the upper end of a distribution and so limit definitive conclusions.

Compared to the other substrates, activities on xyloglucan show an overall increase in average activities. The histogram is also more compressed than those of lichenan and xylan, possibly due to the fact that fully soluble xyloglucan is more easily degraded than the other substrates, so enzymes with activity on xyloglucan are approaching yield saturation. Most relevant to this work, the activities indicate a rank-ordering (Clade 3 > Clade 2 > Clade 1) in median activity that matches the results for lichenan and xylan. Thus it appears that the general catalytic efficiency exhibited by an enzyme in assays of lichenase or xylanase activities transfers to the hydrolysis of xyloglucan. The possibility that open binding clefts in Clades 1 and 2 make these enzymes more generalists with regard to branched substrates remains open, as the tamarind xyloglucan used in this study represents only one possible xyloglucan structure.



**Figure S1.** Per-clade histograms of specific activities on lichenan, xylan, mannan, or xyloglucan. Because the xyloglucan data are not as extensive, activity data for all substrates is presented not corrected for substrate saturation. Histograms are presented in terms of fraction of the tested population.

*Other correlations* — Table S1 reports correlations observed following activity measurements in a second set of assays collected with additional substrates (tamarind xyloglucan, carboxymethylcellulose (CMC), phosphoric acid swollen cellulose (PASC) and beechwood xylan). All reactions in this smaller dataset were carried out for 3 hours at 30 °C, pH 6.

**Table S1.** Correlation Statistics for Additional Substrates and Lichenan or Xylan

Substrate <sup>b</sup>	Enz. <sup>c</sup>	Lichenan <sup>a</sup>				Xylan <sup>a</sup>			
		Pts. <sup>d</sup>	p <sup>e</sup>	$\rho$ <sup>f</sup>	slope <sup>f</sup>	Pts. <sup>d</sup>	p <sup>e</sup>	$\rho$ <sup>f</sup>	slope <sup>f</sup>
Xyloglucan	110	90	< 10 <sup>-3</sup>	0.35 (0.08 – 0.42)	0.0 (0.0 – 0.5)	82	0.7	0.04 (-0.16 – 0.19)	0.0 (-0.3 – 0.3)
CMC	64	47	< 10 <sup>-6</sup>	0.65 (0.51 – 0.76)	0.6 (0.5 – 0.9)	43	0.01	0.38 (0.12 – 0.61)	0.4 (0.0 – 0.9)
PASC	46	20	0.07	0.41 (0.04 – 0.69)	0.5 (0.1 – 1.1)	20	0.6	-0.11 (-0.48 – 0.28)	0.3 (-0.4 – 1.0)
Xylan	64	60	< 10 <sup>-5</sup>	0.56 (0.38 – 0.71)	0.8 (0.6 – 1.0)	53	< 10 <sup>-4</sup>	0.52 (0.30 – 0.70)	0.8 (0.6 – 1.1)

a: Specific activities measured on lichenan and xylan in the main text complete assay (243 enzymes)

b: Substrate assayed in second experimental subset

c: Number of enzymes assayed in second experimental subset

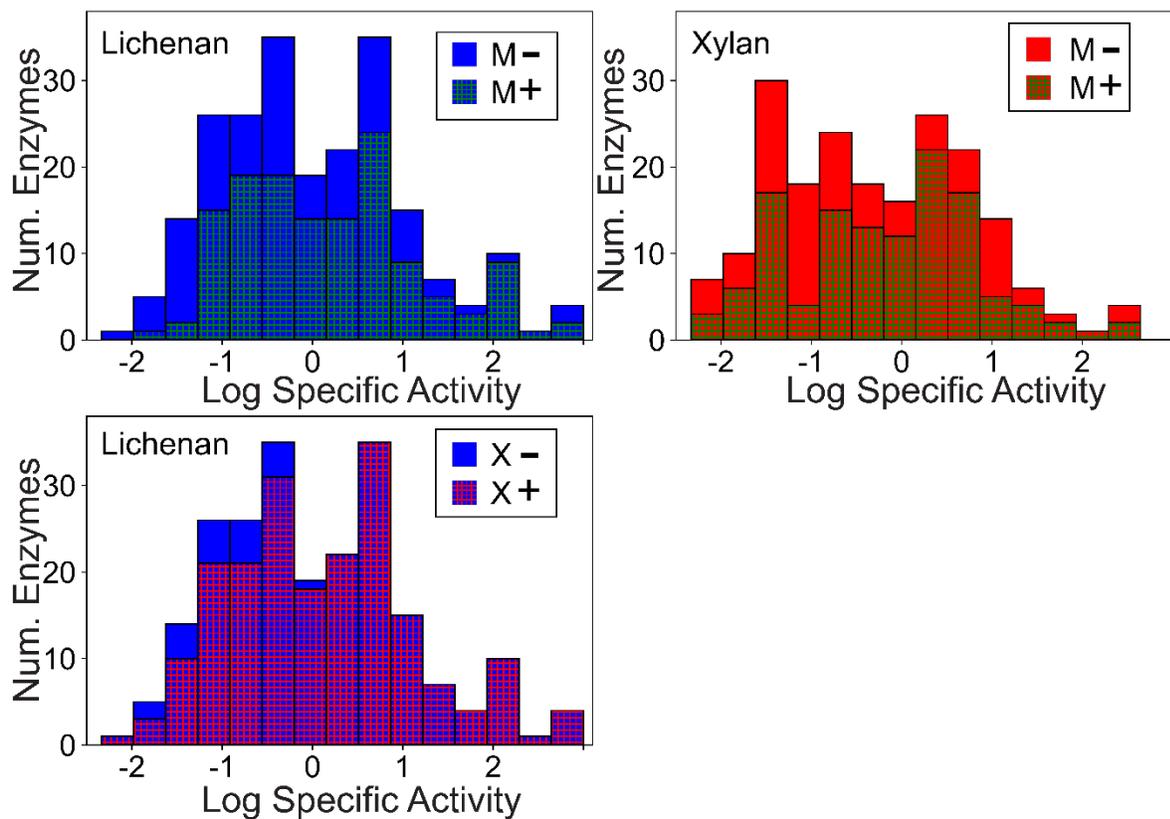
d: Number of datapoints in the log-log comparison dataset (enzymes with both activities > 0)

e: p-value, Spearman  $\rho$ , regression slope

f: Best estimate and 90% confidence interval

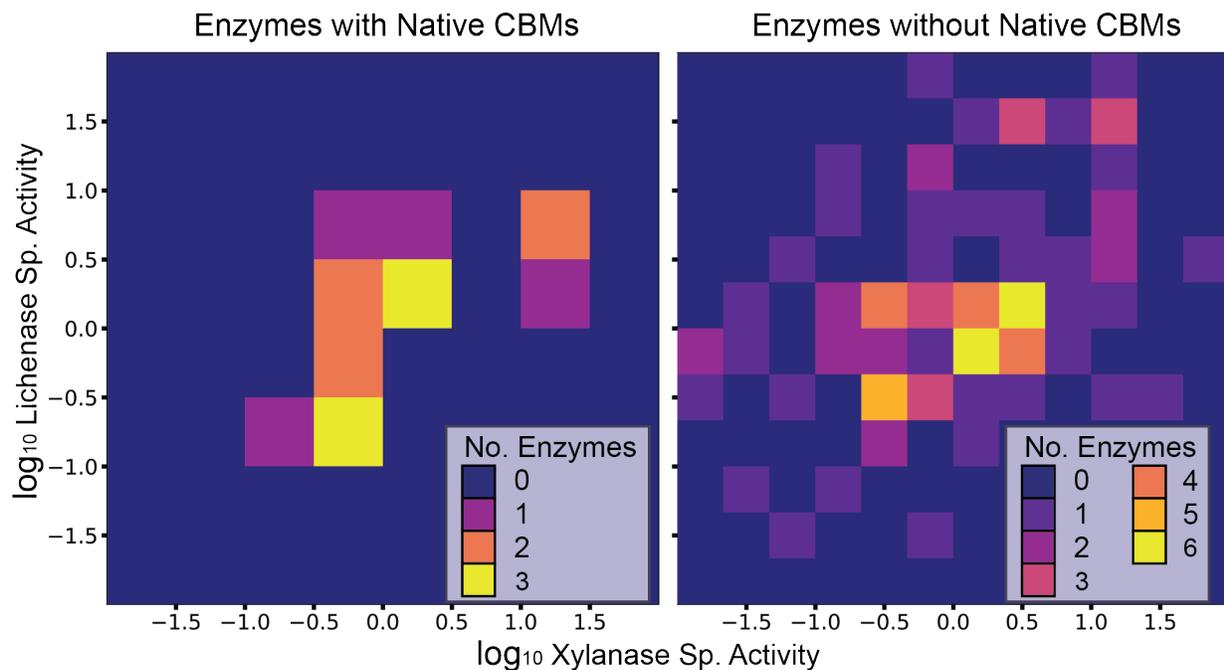
The strongest correlations observed are between specific activities measured on carboxymethylcellulose and lichenan, and between xylan and either lichenan or xylan (the latter serves as an internal validation, since protein translations and substrate preparations were prepared entirely independently of the first set of experiments). Correlations between activities on xyloglucan or PASC and lichenan are also detectable, with the former characterized as rather shallow (log-log slope  $\leq 0.5$ ).

In the full dataset, we do not detect a significant subfamily-wide correlation between activities on mannan and either lichenan or xylan. If there is such a correlation, as mentioned above, both the average activities and histograms suggest it is stronger with lichenan (Figure S2).



**Figure S2.** Histograms of specific activities on lichenan (A and C) or xylan (B), stacked to display frequencies for enzymes with or without mannase activity (M+ vs M-, panels (A) and (B)) or frequencies for enzymes with or without xylanase activity (X+ vs X-, panel (B)).

*Impact of CBM modules* — No obvious functional difference is observed between enzymes possessing or lacking native CBM modules in Clades 2 and 3. Histograms in Figure S3 were generated by merging data from all sub-clades in 2 and 3 following subtraction of the corresponding median activity for that sub-clade. No systematic shift is observable, and a significant difference was not detected (via t-test) between the groups for either xylanase or lichenase activities. While it is possible there are sub-clades or examples in this population where enzyme function is significantly hampered by examining only the catalytic core, the data does not provide strong evidence for this. Moreover, none of the enzymes with no detectable activity were from gene constructs possessing CBM modules (see below).

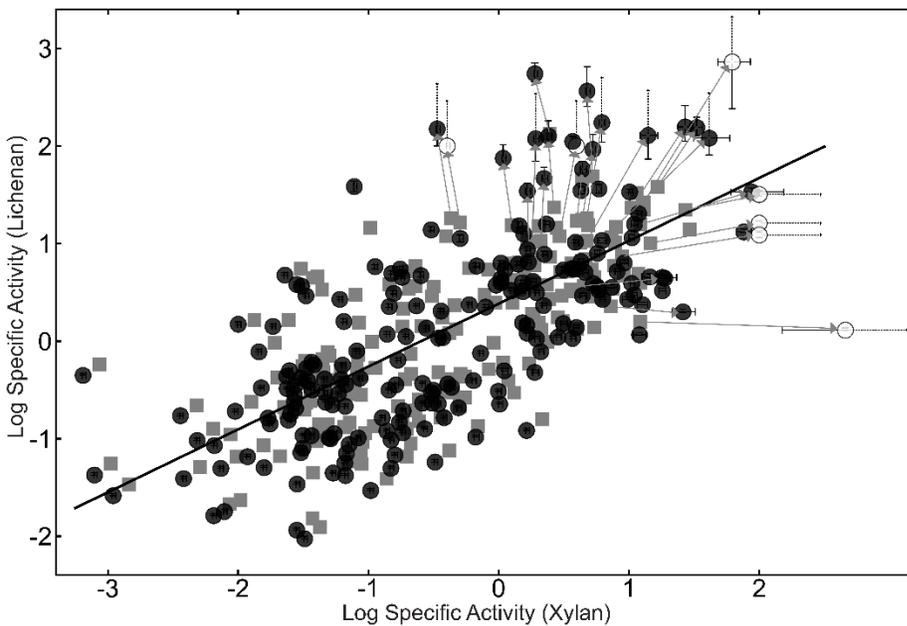


**Figure S3.** 2D histograms plotting the number of Clade 2 and 3 enzymes in a given span of lichenase and xylanase activity, relative to the median activity in the phylogenetic group for each enzyme. For enzymes originating from genes coding for CBMs in addition to the catalytic core (left), bin sizes are 0.5 log units; for enzymes originating from genes without CBMs, bin sizes are 0.33 log units.

In addition to CBMs, genes coding for the enzymes in this study may possess many other types of domains, such as protein localization bacterial immunoglobulin (“Big”) or bacteroidetes-associated carbohydrate-binding often N-terminal (BACON) domains, while other GH5\_4 members harbor dockerin domains for cellulosome association. As mentioned in the main text, several GH5\_4 catalytic domains

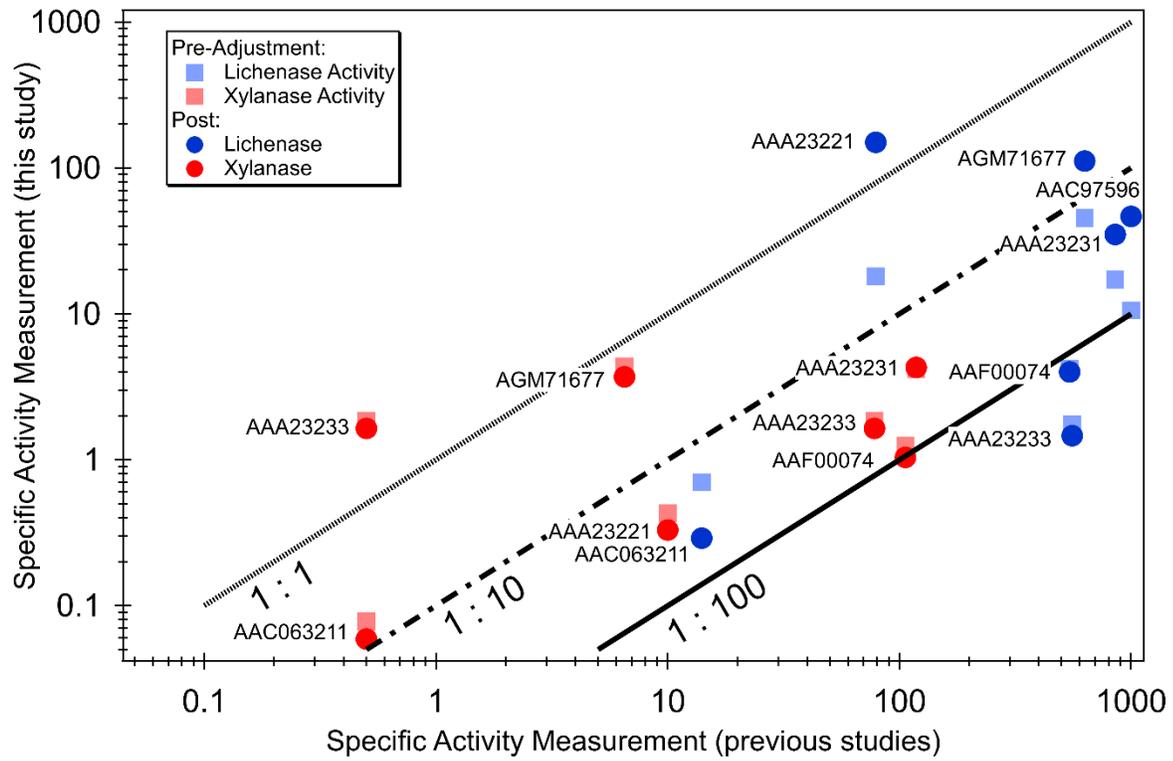
are fused to other enzymatic domains. These can be other GH5\_4 domains (up to four total), a GH26 domain, or a carbohydrate esterase domain, and the functions of these auxiliary enzyme domains may be related plant biomass degradation, such as simultaneous hydrolysis of cellulose and the tightly associated hemicelluloses, or deacetylation of polysaccharide side chains. As was the case with GH5\_4 at the outset of this work, many of these protein families remain mostly under-characterized, and the functions assigned to them are typically based on one or a few examples. A systematic examination of GH5\_4 with the full suite of additional domains observed in nature is beyond the scope of this work but would likely yield tremendous insight into the proteins' true function and utility.

*Compensation for Substrate Depletion and Comparison to Previous Data* — One nuance for the lichenase – xylanase correlation is that some measurements (10% of xylanase data, 22% of lichenase data) were performed after substrates depletions of 10% or more, meaning those data are outside of the linear region typically employed in enzyme kinetics studies (Iakiviak, et al. 2011; McGregor, et al. 2016). Beyond this 10% cutoff, the steady-state approximation becomes invalid and the reaction velocity decreases, so that a single-point specific activity measurement would underestimate the enzyme activity relative to determination in the steady-state regime. As described in Methods, we therefore adjusted enzymes in the main text to a specific activity corresponding to the hypothetical experimental time point corresponding to 5% substrate depletion. This adjustment only significantly impacts data collected after depletions of 10% or more (Figure S4).



**Figure S4.** Lichenase vs xylanase correlation plot pre- and post-adjustment to hypothetical value at 5% substrate depletion. Gray squares, simple specific activity calculations; closed and open circles, adjusted data as plotted in main text Figure 4. Arrows displayed if the plot location is altered by more than 0.5 log units.

Encouragingly, our results also show a consistent relationship to literature measurements throughout the entire dynamic range. Figure S5 displays pre- and post-adjusted data in comparison with previous literature measurements. Measured specific activities on lichenan and xylan are roughly 10 – 100-fold lower in our assays than in several previous studies (Fierobe, et al. 1991; Foong and Doi 1992; Xue, et al. 1992; Liu, et al. 2001; Bianchetti, et al. 2013; McGregor, et al. 2016; Meng, et al. 2017). This is attributed to the lack of shaking in the high-throughput plate assays used here, differences due to scale of reaction, and possibly uncertainties in protein concentration calculations.



**Figure S5.** Comparison between specific activity values in quantitative screen and measurements in previous studies

Complete phylograms for each main group with detailed annotation are displayed below in Figure S6A – D (figures corresponding to related subfamilies, and Clades 1, 2, and 3). Bayesian MCMC probabilities ('supports') are listed above each node. Activity annotations from this study are included to the right of each node tip as in Figure 4 of the main text, with an added purple square similarly denoting that xyloglucanase activity was measured. As for the other activities, the size of the xyloglucanase symbol corresponds to its activity lying within the first, second or third quantile. EC numbers listed in the CAZy database are also listed. In the informational table at the right-hand side of the figure, the accession code and Pfam motif annotation for each gene are listed. Motifs are plotted to highlight features most relevant to this work: purple rectangles, GH5; yellow rectangles, other GH families; purple triangles, CBM X2 motif (a.k.a. CBM46 and the related Ig-like domain); tan triangles, other CBM motifs; brown ellipses, all other motifs. Locations of crystal structures and representative PDB accession codes are also marked between the relevant node tip and table entry.

Letters in the final column list isolation notes as provided in NCBI database: W= aquatic sample, S=terrestrial sample; R= rumen or other host-associated sample. Aside from handful of examples, Clade 1 enzymes have been found in soil or aquatic samples. Clades 2 and 3 are heavily represented by rumen and gut symbionts; of evolutionary note, even within subclades 2C, 2D and 3B, which are dominated by free living organisms, the enzymes that are most closely related to the common ancestor come from host-associated organisms.

Closely-related subfamilies and study outgroup representatives from subfamily 2 (AAC19169.1) and subfamily 5 (AFY52522.1) were included in the experimental set and the tree-building process. In Region I, there is a representative crystal structure corresponding to the protein with Genbank accession code ABN52701.1. This sequence was not included in the tree-building process, but it is most closely related to ACL75216.1 and AKP16689.1. It is marked as an added hypothetical node – table connecting line.

Fig. S6A

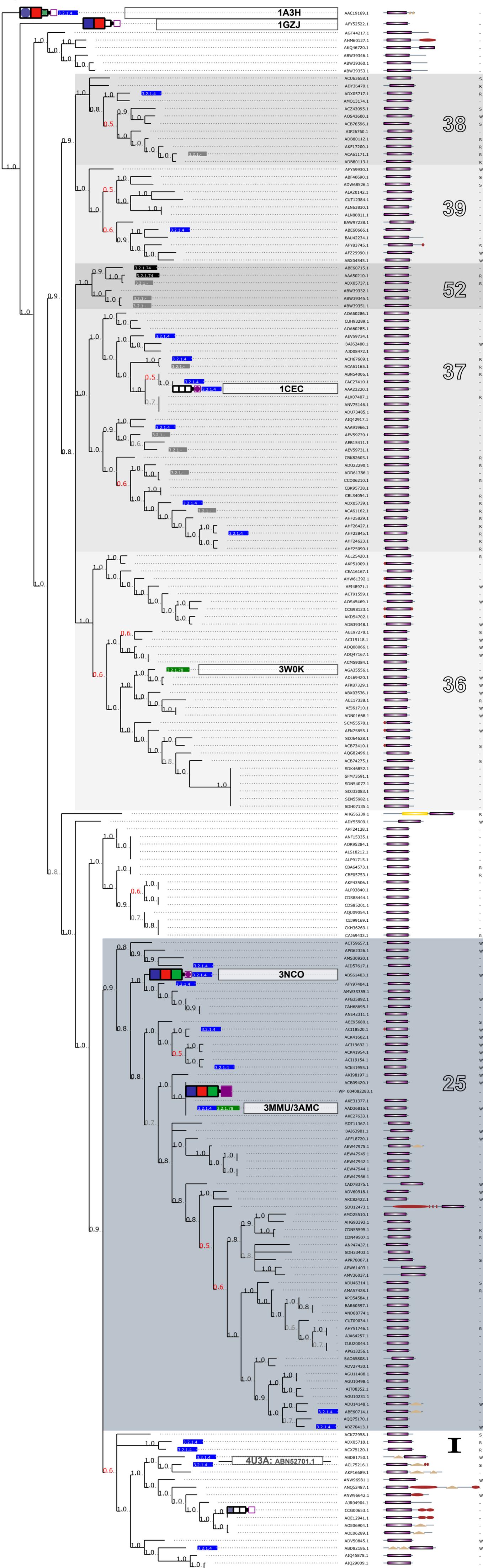


Fig. S6B

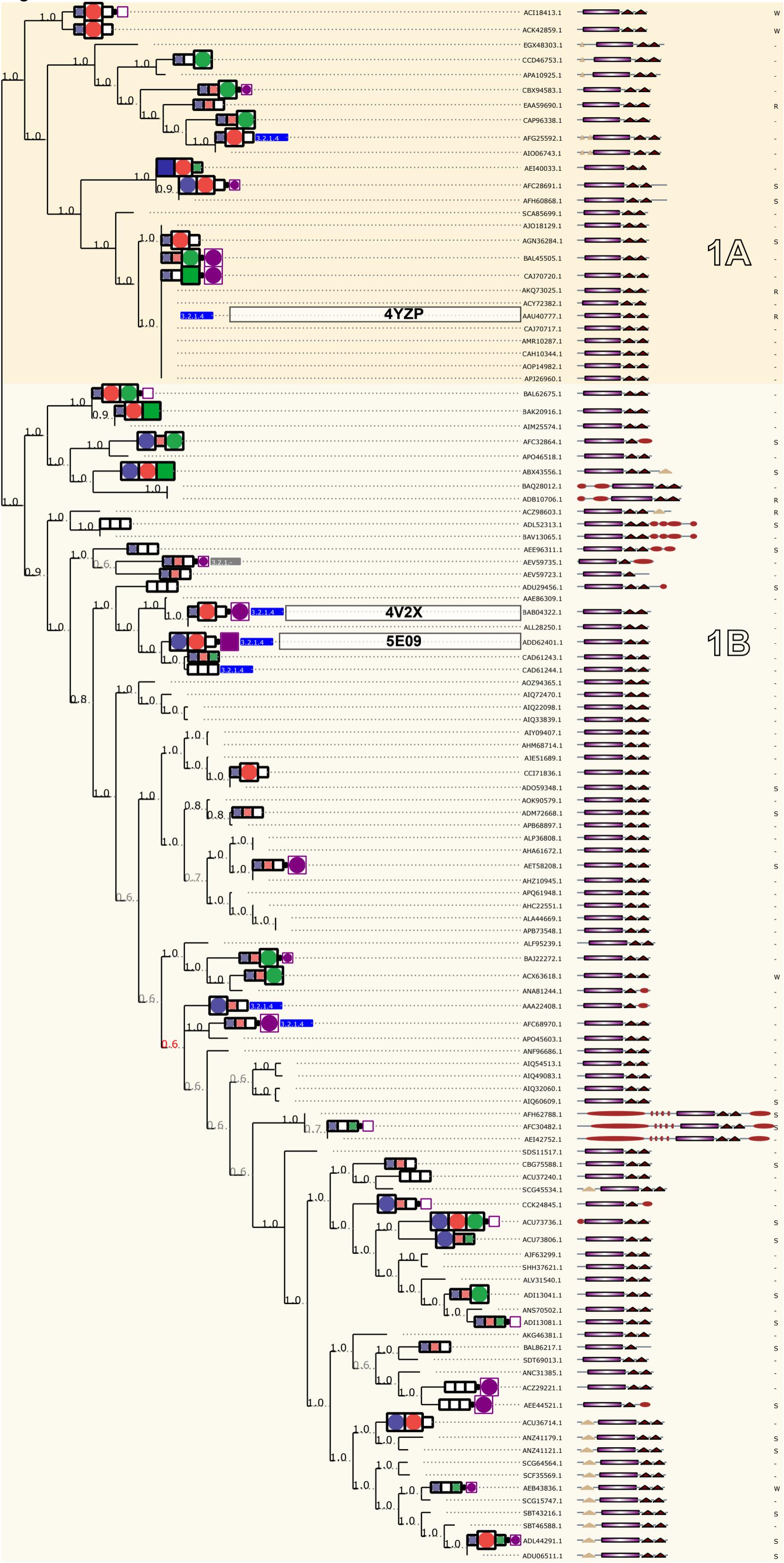


Fig. S6C

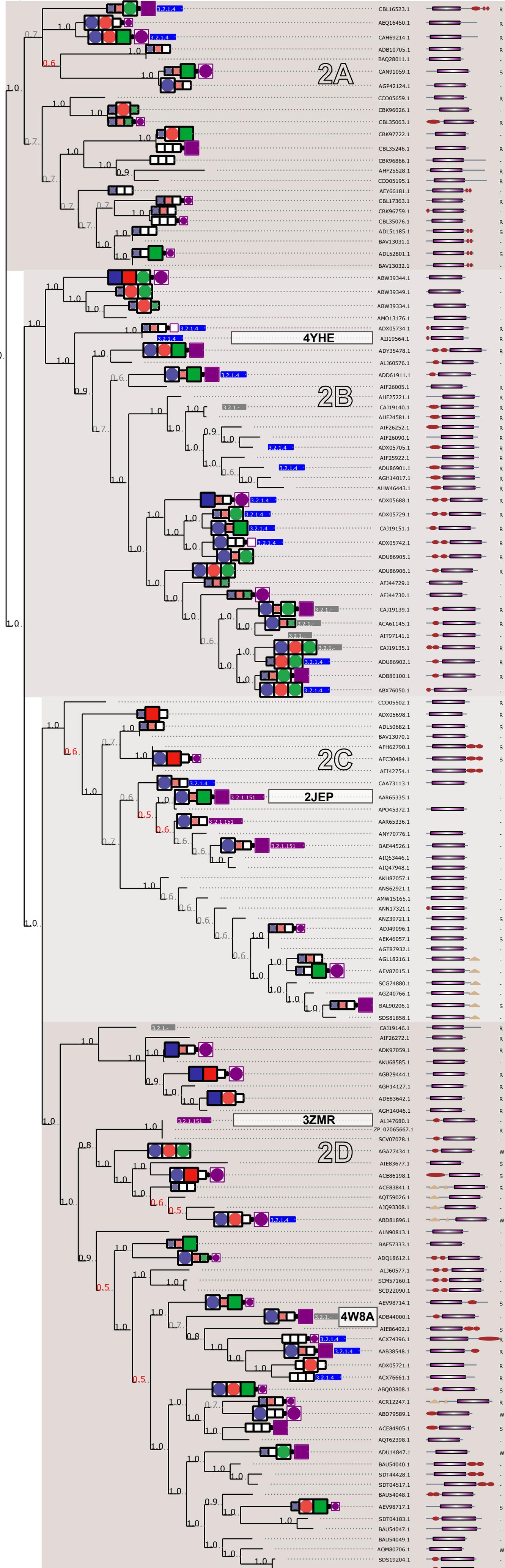
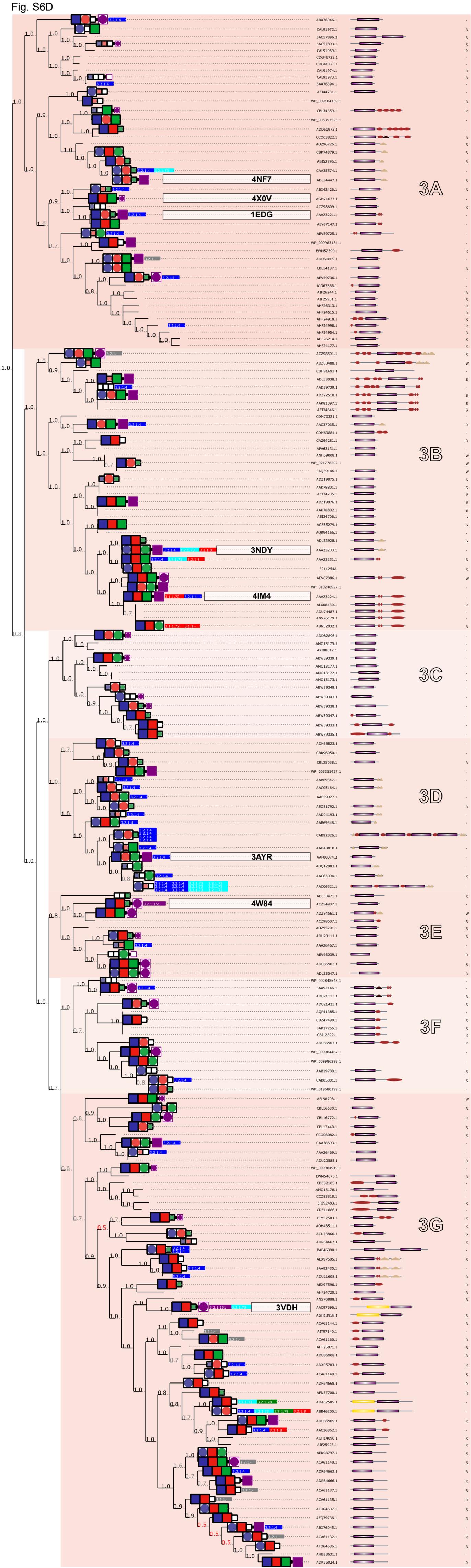


Fig. S6D



## References

- Bianchetti CM, Brumm P, Smith RW, Dyer K, Hura GL, Rutkoski TJ, Phillips GN, Jr. 2013. Structure, dynamics, and specificity of endoglucanase D from *Clostridium cellulovorans*. *J Mol Biol* 425:4267-4285.
- Fierobe HP, Gaudin C, Belaich A, Loutfi M, Faure E, Bagnara C, Baty D, Belaich JP. 1991. Characterization of endoglucanase A from *Clostridium cellulolyticum*. *J Bacteriol* 173:7956-7962.
- Foong FCF, Doi RH. 1992. Characterization and Comparison of *Clostridium-Cellulovorans* Endoglucanases-Xylanases Engb and Engd Hyperexpressed in *Escherichia-Coli*. *Journal of Bacteriology* 174:1403-1409.
- Iakiviak M, Mackie RI, Cann IK. 2011. Functional analyses of multiple lichenin-degrading enzymes from the rumen bacterium *Ruminococcus albus* 8. *Appl Environ Microbiol* 77:7541-7550.
- Liu J, Tsai C, Liu J, Cheng K, Cheng C. 2001. The catalytic domain of a *Piromyces rhizinflata* cellulase expressed in *Escherichia coli* was stabilized by the linker peptide of the enzyme. *Enzyme Microb Technol* 28:582-589.
- McGregor N, Morar M, Fenger TH, Stogios P, Lenfant N, Yin V, Xu X, Evdokimova E, Cui H, Henrissat B, et al. 2016. Structure-Function Analysis of a Mixed-linkage beta-Glucanase/Xyloglucanase from the Key Ruminant Bacteroidetes *Prevotella bryantii* B(1)4. *J Biol Chem* 291:1175-1197.
- Meng DD, Liu X, Dong S, Wang YF, Ma XQ, Zhou H, Wang X, Yao LS, Feng Y, Li FL. 2017. Structural insights into the substrate specificity of a glycoside hydrolase family 5 lichenase from *Caldicellulosiruptor* sp. F32. *Biochem J* 474:3373-3389.
- Xue GP, Gobius KS, Orpin CG. 1992. A Novel Polysaccharide Hydrolase Cdna (Celd) from *Neocallimastix-Patriciarum* Encoding 3 Multifunctional Catalytic Domains with High Endoglucanase, Cellobiohydrolase and Xylanase Activities. *Journal of General Microbiology* 138:2397-2403.