

# Genetic and functional diversity of $\beta$ -N-acetylgalactosamine residue-targeting glycosidases expanded by deep-sea metagenome

## Authors:

Tomomi Sumida<sup>1\*</sup>, Satoshi Hiraoka<sup>1</sup>, Keiko Usui<sup>1</sup>, Akihiro Ishiwata<sup>2</sup>, Toru Sengoku<sup>3</sup>, Keith A Stubbs<sup>4</sup>, Katsunori Tanaka<sup>2,5</sup>, Shigeru Deguchi<sup>1</sup>, Shinya Fushinobu<sup>6\*</sup>, Takuro Nunoura<sup>1</sup>.

## Affiliations:

1. Research Center for Bioscience and Nanoscience, Research Institute for Marine Resources Utilization, Japan Agency for Marine-Earth Science and Technology (JAMSTEC), Yokosuka, Japan.

2. RIKEN, Cluster for Pioneering Research, Wako, Saitama, Japan.

3. Department of Biochemistry, Yokohama City University Graduate School of Medicine, Kanazawa-ku, Yokohama, Japan.

4. School of Molecular Sciences, The University of Western Australia, Crawley, WA, Australia.

5. Department of Chemical Science and Engineering, Tokyo Institute of Technology, Meguro-ku, Tokyo, Japan.

6. Graduate School of Agricultural and Life Sciences, The University of Tokyo, Bunkyo-ku, Tokyo, Japan.

## \*Corresponding authors:

Tomomi Sumida, Email:sumidat@jamstec.go.jp,

Shinya Fushinobu, Email:asfushi@mail.ecc.u-tokyo.ac.jp

## Abstract

$\beta$ -*N*-Acetylgalactosamine-containing glycans play essential roles in several biological processes, including cell adhesion, signal transduction, and immune responses.  $\beta$ -*N*-Acetylgalactosaminidases hydrolyze  $\beta$ -*N*-acetylgalactosamine linkages of various glycoconjugates. However, their biological significance remains ambiguous, primarily because only one type of enzyme, exo- $\beta$ -*N*-acetylgalactosaminidases that specifically act on  $\beta$ -*N*-acetylgalactosamine residues, has been documented so far. In this study, we identified three novel glycoside hydrolase families distributed among all three domains of life and characterized eight novel  $\beta$ -*N*-acetylgalactosaminidases and  $\beta$ -*N*-acetylhexosaminidase through sequence-based screening of deep-sea metagenomes and subsequent searching of public protein databases. Despite low sequence similarity, the crystal structures of these enzymes demonstrate that all enzymes share a prototype structure and diversify their substrate specificities (endo-, dual-endo/exo-, and exo-) through the accumulation of mutations and insertional amino acid sequences. The diverse  $\beta$ -*N*-acetylgalactosaminidases reported in this study could facilitate the comprehension of their structures and functions and present novel evolutionary pathways for expanding their substrate specificity.

Beta-*N*-acetylgalactosamine ( $\beta$ -GalNAc)-containing glycans, such as glycoconjugates of polysaccharides (1), glycolipids (2, 3), *N*- and *O*-linked glycans (4, 5), *O*-antigen (6), and chondroitin sulfate (7), are ubiquitous and crucially contribute to various biological processes, including cell adhesion, signal transduction, cross-interactions with functional membrane components, formation of the cell envelope and maintenance of its stability, immunomodulation, and immune responses (1–7). The regulatory function of these glycans is attributed to their structural diversity, which differ in carbohydrate constituents (namely, glucose, galactose) and molecular architecture ( $\alpha$ - or  $\beta$ -linkages, linear or branched) (1–8).

Beta-*N*-acetylgalactosaminidases ( $\beta$ -NGAs) hydrolyze the different  $\beta$ -GalNAc linkages of various glycans to modulate the length, combination, and abundance of glycans. This catalytic activity requires the  $\beta$ -NGAs to possess diverse substrate specificities, but only two enzymes, *exo*- $\beta$ -NGA and *exo*- $\beta$ -*N*-acetylhexosaminidase (*exo*- $\beta$ -HEX), possessing distinctive substrate specificity and sequence, have been evidenced to hydrolyze  $\beta$ -GalNAc. *Exo*- $\beta$ -NGAs have a strict substrate specificity for the non-reducing terminal  $\beta$ -GalNAc and are classified into the glycoside hydrolase (GH) family GH123 (9) of the Carbohydrate-Active Enzymes (CAZy) database (10, 11). As  $\beta$ -GalNAc is prevalent in various glycans across the three domains of life (Bacteria, Archaea, and Eukarya) and in different ecological niches (1–8),  $\beta$ -NGA activity is expected to follow a similar distribution pattern. However, these enzymes have been identified in only three bacterial species, specifically associated with microbe-host interactions in both terrestrial soil and human gut environments (namely, NgaP from *Paenibacillus* sp. TS12, CpNga123 from *Clostridium perfringens*, and BvGH123 from *Phocaeicola vulgatus*) (9, 12, 13), with no reported origin in archaea or eukaryotes, and *endo*- $\beta$ -NGA has not been reported. Meanwhile, *exo*- $\beta$ -HEXs hydrolyze the non-reducing terminal of  $\beta$ -GalNAc as well as  $\beta$ -GlcNAc, and are classified into the family GH20 (14–16). More than 100 *exo*- $\beta$ -HEXs from all three domains of life have been functionally analyzed, with no report of *endo*- $\beta$ -HEXs.

Given the limited extant knowledge on  $\beta$ -NGAs, it is imperative to further identify and functionally characterize novel  $\beta$ -NGAs. These endeavors are critical for comprehensively understanding the complex phenomena associated with  $\beta$ -GalNAc-mediated biological processes. Recently, culture-independent metagenomic exploration of novel glycosidases has substantially augmented our conception of carbohydrate-related enzymes. Function-based screening of diverse biological resources revealed several novel glycosidase families, including GH148 from volcanic soil (17), GH156 from a thermal hot spring (18), GH165

from agricultural soil (19), and GH173 and CBM89 from the capybara intestine (20). Moreover, sequence-based screening has enabled the analysis of much larger metagenomic sequencing datasets than functional screening, yielding more candidate sequences and facilitating the discovery of novel enzymes with distinct characteristics (21).

The deep-sea environments, characterized by unique features and distinct bacterial flora (22), have rarely been surveyed owing to the challenges associated with sampling from these regions. The deep-sea metagenome is a promising frontier for enzyme discovery (23, 24). Therefore, here we aimed to use a sequence-based screening approach of deep-sea microbial assemblages to explore the functional diversity of  $\beta$ -NGA activity. Our deep-sea sediment metagenomic (DSSM) analysis and domain search yielded three novel  $\beta$ -NGA gene families that are phylogenetically distinct from GH123 exo- $\beta$ -NGAs. The biochemical and structural characterization of these enzymes not only unveiled their functional diversity but also shed light on their monophyletic evolutionary history, providing valuable insights into the mechanisms underlying  $\beta$ -GalNAc-mediated biological processes.

## Results

### Discovery of a novel $\beta$ -NGA with both endo- and exo-glycosidase activities

Using four metagenomic datasets derived from microbial assemblages in deep-sea abyssal sediments and a domain-based search, we retrieved three candidate complete coding sequences (CDSs) (Gene ID: *dssm\_1*–3, tentative Protein ID: DSSM\_1–3), which exhibited low sequence identity (15–26%) to all known GH123 exo- $\beta$ -NGA genes (*NgaP*, *CpNga123*, and *BvGH123*) (Extended Data Fig. 1a, Supplementary Data S1). Sequence alignments demonstrated that the consecutive catalytic “DE” motif of family GH123, comprising an aspartic acid (stabilizer of the 2-acetamido group of the substrate) and a glutamic acid (acid/base) (9), was present in the *dssm\_2* and *dssm\_3* sequences but not in *dssm\_1* (Extended Data Fig. 1b, green box and asterisk).

AlphaFold2 (25) was used to predict the structure of the candidate proteins (Fig. 1a). DSSM\_1 was structurally distinct from the GH123 exo- $\beta$ -NGAs (Fig. 1a, b) and consisted of five domains, among which only domain 2 ( $\beta$ -sandwich) displayed a degree of structural similarity to the N-terminal domain of GH123 exo- $\beta$ -NGAs. Domain 4 ( $(\beta/\alpha)_8$  barrel) was similar to cycloisomaltooligosaccharide glucanotransferase (PDB, 3WNM). By contrast, the predicted structures of DSSM\_2 and DSSM\_3 were similar to those of GH123- $\beta$ -NGAs, and

they all shared DUF4091, a presumed domain region whose function remains undetermined (Fig. 1b, green, Extended Data Fig. 1b, underlined).

Next, the  $\beta$ -NGA candidate sequences *dssm\_2* and *dssm\_3* were selected and heterologously expressed in *Escherichia coli*. The encoding sequences lacking the predicted signal peptides (Supplementary Data S1, bold and underlined) were cloned into an expression vector. Although an expression construct of *dssm\_2* failed to yield a soluble protein, the enzyme encoded by *dssm\_3* was solubilized, and a purified protein was successfully obtained (renamed Protein ID: NgaDssm). Assays using various *p*NP-substrates revealed that NgaDssm was active on GalNAc- $\beta$ -*p*NP, but not on GlcNAc- $\beta$ -*p*NP or GalNAc- $\alpha$ -*p*NP, indicating that the enzyme possessed exo- $\beta$ -NGA but not exo- $\beta$ -HEX activity. Furthermore, the enzyme hydrolyzed Gal $\beta$ 1-3GalNAc- $\beta$ -*p*NP but not Gal- $\beta$ -*p*NP, demonstrating an additional disaccharide-releasing endo- $\beta$ -NGA activity (Table 1). These findings suggest that NgaDssm is a novel endo/exo- $\beta$ -NGA.

### Phylogenetic diversity of $\beta$ -NGAs

The uncharacterized domain DUF4091 was conserved in the GH123 exo- $\beta$ -NGAs and a novel endo/exo-NGA gene sequence. Consequently, we utilized DUF4091 as a query to further identify  $\beta$ -NGA genes. We retrieved 734 sequences containing DUF4091 from the Pfam protein family database. The catalytic DE motif was highly conserved (94%), and thus, majority of these genes were expected to encode enzymes possessing  $\beta$ -NGA activity. A phylogenetic tree of these sequences, along with the three known GH123 exo- $\beta$ -NGA genes and the two deep-sea  $\beta$ -NGA candidates (Supplementary Data S2), identified four major groups: GH123, Group 1, Group 2, and Group 3 (Fig. 1c). Group 3 was further segregated into two subgroups (Group 3-1 and Group 3-2) based on the phylogenetic analysis, with NgaDssm belonging to Group 3-1 (Fig. 1c, d).

To examine the  $\beta$ -NGA sequence alignment and activity distributed among all three domains of life, 14 representative sequences of each novel group were selected from model plants, archaea, and various bacterial phyla, as follows: for Group 1, *Cohnella abietic* from Bacillota (NgaCa), *Meiothermus granaticius* from Deinococcota (NgaMg), *Nostoc punctiforme* (NgaNp), *Clyndrospermum stagnale* (NgaCs), and *Stanieria cyanosphaera* (NgaSc) from Cyanobacteria; for Group 2, *Arabidopsis thaliana* (NgaAt), *Glycine max* (NgaGm), and *Oryza sativa* (NgaOs) from plants; for Group 3-1; *Candidatus Bathyarchaea archaeon* B24 from Thermoproteota in Archaea from the hydrothermal vent microbiome

(NgaBa) and *Chitinophaga pinensis* from Bacteroidota (NgaCp); for Group 3-2; *Brachybacterium faecium* (NgaBf) and *Bifidobacterium longum* subsp. *infantis* (NgaBl) from Actinomycetota, *Lacticaseibacillus yichunensis* (NgaLy), and *Paenibacillus* sp. TS12 (NgaP2) from Bacillota (Supplementary Data S2 and 3). Overall, the sequence identities between each group of 14 candidate genes and the three GH123 genes were low (12–26%) based on sequence alignment (Extended Data Fig. 2a), and only nine amino acids were entirely conserved (Fig. 1d, red background). DUF4091 was a comparatively well-conserved region, comprising four of the nine strictly conserved amino acids, with its sequence identity being at least 10% higher than that of the full-length sequence (Extended Data Fig. 2b, c).

### Substrate specificity of $\beta$ -NGA candidates

Recombinant expression vectors for the new candidate genes (except for NgaGm and NgaOs) were constructed to assess the activity of the identified enzymes (Supplementary Data S3). Although four expression constructs (NgaNp, NgaCs, NgaSc, and NgaCp) failed to yield soluble proteins, the remaining eight (NgaCa, NgaMg, NgaAt, NgaBa, NgaBf, NgaBl, NgaLy, and NgaP2) successfully expressed soluble enzymes (Extended Data Fig. 2a, red letter) and were subjected to protein purification for subsequent enzymatic assays (Extended Data Fig. 3a, Table 1).

The substrate specificity of the expressed enzymes was assayed using various synthetic substrates (Table 1). Surprisingly, Group 1 NgaCa displayed strictly endo-type enzymatic activity and acted solely on Gal $\beta$ 1-3GalNAc- $\beta$ -pNP. The other Group 1 protein, NgaMg, demonstrated no activity against any of the tested substrates.

In Group 2, NgaAt was active against GalNAc- $\beta$ -pNP and GlcNAc- $\beta$ -pNP, but not against Gal $\beta$ 1-3GalNAc- $\beta$ -pNP, indicating that it possessed exo- $\beta$ -HEX activity. However, NgaAt did not share notable sequence similarity with any of the three types of exo- $\beta$ -HEXs (HEXO1–3) classified into the GH20 family from *A. thaliana* (Extended Data Fig. 4a, b) (26, 27). Thus, NgaAt is a novel family of exo- $\beta$ -HEX.

NgaBa in Group 3-1 was active on both GalNAc- $\beta$ -pNP and Gal $\beta$ 1-3GalNAc- $\beta$ -pNP, similar to NgaDssm, suggesting that members of Group 3-1 share endo/exo-type  $\beta$ -NGA activity. In Group 3-2, NgaBf, NgaBl, and NgaLy were active only on GalNAc- $\beta$ -pNP, highlighting their strict exo- $\beta$ -NGA functionality. NgaP2 was active against GalNAc- $\beta$ -pNP and had weak activity against GlcNAc- $\beta$ -pNP. Collectively, these results showed that the enzymes in each group exhibit a characteristic endo- and/or exo-type cleavage mode and



substrate specificity.

NgaP from GH123 is an excellent tool for detecting sulfatase deficiency, as it does not act on substrates sulfated at positions 4 or 6 of GalNAc (28), and has been successfully used for screening of mucopolysaccharidosis (a metabolic disorder caused by the accumulation of mucopolysaccharides) in newborns (28, 29). We evaluated whether the enzymes identified herein possessed GalNAc4S or GalNAc6S cleavage activity. Using GalNAc4S- $\beta$ -4MU and GalNAc6S- $\beta$ -4MU (Table 1), we observed that NgaDssm, NgaAt, and NgaP2 acted on GalNAc4S- $\beta$ -4MU. In particular, NgaDssm displayed approximately two-fold higher activity against GalNAc4S- $\beta$ -4MU than against GalNAc- $\beta$ -4MU.

### Generic properties of the novel $\beta$ -NGAs

Next, we evaluated the enzyme characteristics (Table 1, Extended Data Fig. 3). The optimal pH for the eight novel enzymes was in the range 5.0–7.0. NgaBa depicted a particularly broad optimal pH range, with relative activity above 90% at pH 5.5–7.5. Metal ions generally did not affect enzyme activity, except that of NgaLy, which was inhibited by  $\text{Ni}^{2+}$ ,  $\text{Co}^{2+}$ , and  $\text{Zn}^{2+}$  (Table 1, Extended Data Fig. 3b, c). Overall, the enzymes exhibited moderate temperature optima (25–45°C) except for NgaBa, an enzyme from a hydrothermal vent microbiome (30), which displayed an extremely high optimal temperature (70°C) and  $T_m$  value (93.1°C) (Extended Data Fig. 3d, e). Thus, NgaBa is the first thermostable  $\beta$ -NGA reported to date. The  $K_m$  and  $k_{cat}$  of each enzyme for the most preferred substrate (GalNAc- $\beta$ -pNP or Gal $\beta$ 1-3GalNAc- $\beta$ -pNP) was also examined (Extended Data Fig. 3f), wherein NgaLy possessed the highest  $k_{cat}/K_m$  value.

### Substrate specificity for oligosaccharides

The substrate specificity of each enzyme was explored using various oligosaccharides as substrates (Extended Data Fig. 5a). NgaCa and NgaMg from Group 1 and NgaDssm from Group 3-1 surprisingly did not act on GA1 and Gb5 oligosaccharides, although these two oligosaccharides shared the non-reducing end structure Gal $\beta$ 1-3GalNAc- $\beta$ - with Gal $\beta$ 1-3GalNAc- $\beta$ -pNP (Extended Data Fig. 5b, d). NgaAt acted on GalNAc $\beta$ 1-3Gal but not on GalNAc $\beta$ 1-4Gal (Extended Data Fig. 5c). By contrast, NgaBa was functional against both GA1 and Gb5 (Extended Data Fig. 5d, right lanes 2 and 5), displaying a more rapid digestion of GA1 over that of GA2 (Extended Data Fig. 5d, right lanes 2 and 4). The endo-activity of NgaBa against  $\beta$ -GalNAc located inside the oligosaccharide was stronger than its exo-

activity against  $\beta$ -GalNAc at the non-reducing ends, but this preference was reversed when *p*NP-substrates were used. Moreover, NgaBa more effectively degraded Gb5 than Gb4 (endo- > exo-activity) and displayed enzymatic activity against Gal $\beta$ 1-3GalNAc- $\beta$ - and GalNAc $\alpha$ 1-3GalNAc- $\beta$ - (Extended Data Fig. 5d). The NgaBf, NgaBl, NgaLy, and NgaP2 activities were similar to those of GH123 exo- $\beta$ -NGAs, as they broke down only linear oligosaccharides with  $\beta$ -GalNAc at the non-reducing terminus (Extended Data Fig. 5e). The results for NgaBf and NgaLy indicated that these enzymes preferred GalNAc $\beta$ 1-4Gal to GalNAc $\beta$ 1-3Gal. As only a few  $\beta$ -GalNAc-containing oligosaccharides are commercially available, the natural substrates of NgaCa, NgaMg, and NgaDssm were not identified.

### X-ray crystal structures of the $\beta$ -NGA from novel groups

Although the novel  $\beta$ -NGAs and GH123 exo- $\beta$ -NGAs were divided into four major groups based on the phylogenetic analysis, their predicted overall structure was similar. Based on substrate preferences, several of these enzymes considerably differed from GH123 exo- $\beta$ -NGAs in terms of substrate specificity. To understand the relationship between structural variation and substrate specificity, crystallization screening was performed on all novel  $\beta$ -NGAs, and X-ray crystallographic analyses of apo- and/or ligand-bound forms were successfully performed on the following enzymes (NgaCa [Group 1], NgaAt [Group 2], NgaDssm [Group 3-1], NgaLy, and NgaP2 [Group 3-2]; nine forms of five enzymes in total) (Fig. 2 and Supplementary Tables 1–5). The overall structure of these enzymes is similar to that of CpNga123 and BvGH123, which consist of a  $\beta$ -sandwich domain at the N-terminus and a ( $\beta/\alpha$ )<sub>8</sub>-barrel domain encompassing the catalytic region (Fig. 2a, Supplementary Figs. 1–5). The structure of NgaDssm (Group 3-1) bore striking resemblance to that of NgaLy and NgaP2 (Group 3-2) (root-mean-square distance of C $\alpha$  atoms [rmsd] = 1.5–2.3 Å), while the structures of NgaLy and NgaP2 are identical (rmsd = 1.0 Å) (Extended Data Fig. 6a). The additional N-terminal domain consisting of approximately 80 amino acids, found only in Group 2 enzymes, was not modeled in the crystal structure of NgaAt due to disorder (Supplementary Fig. 2).

Among Groups 1, 2, and 3 and GH123, the DE motif is located at the same position, suggesting that these enzymes adopt the substrate-assisted catalysis reported for GH123 exo- $\beta$ -NGAs (Fig. 2b). The conserved DUF4091 domain in these enzyme groups is located at the innermost position of the enzyme between the N-terminal domain and the ( $\beta/\alpha$ )<sub>8</sub>-barrel domain (Fig. 2b, purple). This suggested that DUF4091 most likely plays a role in defining



the location of N-terminal and ( $\beta/\alpha$ )<sub>8</sub>-barrel domains as a “central pillar”. Intriguingly, nine amino acid residues conserved among the groups (Fig. 1d, red background) are found at the same positions in the structure (Fig. 2c). Among the conserved amino acids, D95, W115, Y466, and R478 (amino acid numbers of NgaLy [Extended Data Fig. 6b]) are located at the interface between the N-terminal domain and DUF4091 (Fig. 2b, right). D95 and W115, located at the N-terminal domain, fit into the surface pocket of DUF4091, and the side chain of D95 formed hydrogen bonds with Y466 and R478 located at DUF4091 (Extended Data Fig. 6c). Therefore, D95, W115, Y466, and R478 are potentially important for maintaining structural integrity as they link the N-terminal and ( $\beta/\alpha$ )<sub>8</sub>-barrel domains (Fig. 2b, blue letter). G427 is positioned at the entry point for the eighth  $\beta$ -sheet into the ( $\beta/\alpha$ )<sub>8</sub> barrel domain, with the space barely the size of the glycine residue (Extended Data Fig. 6d). The DE motif (D327 and E328), Y392, and W431 are located around the substrate and are involved in substrate recognition (Fig. 2b, red letter). Previous studies have reported that point mutations in the DE motif decrease catalytic activity (9, 12). Therefore, we examined the remaining seven conserved amino acids and constructed corresponding alanine mutants (Extended Data Fig. 6e). The point mutants D95A, W115A, Y466A, and R478A were expressed in *E. coli*, but all proteins precipitated. For G427, a Val mutant was also constructed based on structural information that suggested the conversion of Gly to Ala would be tolerated (Extended Data Fig. 6d). In this case, a very small amount of solubilized enzyme was obtained for G427A but not for G427V (Extended Data Fig. 6e, blue letter). The Y392A and W431A mutants were purified as soluble enzymes (Extended Data Fig. 6e, red letter), but their catalytic activity was considerably reduced (Extended Data Fig. 6f). These findings imply that amino acid substitutions severely impact the structural stability (D95, W115, Y466, R478, and G427) and substrate recognition (Y392 and W431) of the enzyme.

## Essential structural elements for endo- and exo-specificity

We compared the overall structure and substrate-binding sites among the groups to identify the structural elements governing endo- and exo-specificity of these enzymes (Fig. 3). NgaCa in Group 1 had the simplest structure among the enzymes analyzed (Fig. 3a, Supplementary Fig. 1). The active site of NgaCa possesses a cleft shape that allows oligosaccharides to pass through to the -2 subsite, enabling oligosaccharide binding and endo-type  $\beta$ -NGA activity. This structural feature indicated that Group 1 enzymes act on the inner regions of longer sugar chains in addition to disaccharides from the non-reducing end

(e.g., Gal $\beta$ 1-3GalNAc- $\beta$ -pNP).

The active site of NgaAt in Group 2 also possesses a cleft shape (Fig. 3b), and the second  $\beta$ -sheet of the ( $\beta/\alpha$ )<sub>8</sub> barrel is longer than that found in the other enzymes, followed by an additional  $\beta$ -sheet (Fig. 3b and Supplementary Fig. 2, orange). Consequently, the space around the 3-OH group of  $\beta$ -GalNAc in the -2 subsite is narrower than that in Group 1. Therefore, NgaAt was presumed to act on Gal $\beta$ 1-4GalNAc $\beta$ - but not on Gal $\beta$ 1-3GalNAc $\beta$ -. Since Gal $\beta$ 1-4GalNAc- $\beta$ -pNP was commercially unavailable, Gal $\beta$ 1-4GlcNAc- $\beta$ -pNP was used as an alternative substrate to test this hypothesis, and we observed that Gal $\beta$ 1-4GlcNAc- $\beta$ -pNP, but not Gal $\beta$ 1-3GalNAc- $\beta$ -pNP, was degraded (Extended Data Fig. 4c). Thus, NgaAt demonstrates more diverse activity than expected and possesses endo/exo- $\beta$ -HEX functionality, making it the first  $\beta$ -HEX with this characteristic.

The active site of the endo/exo-type NgaDssm from Group 3-1 has a cleft structure that enables passage to the -2 subsites, as in NgaCa, but it also possesses two additional  $\beta$ -hairpins from the second and third  $\beta$ -sheets of the ( $\beta/\alpha$ )<sub>8</sub> barrel domain (Fig. 3c and Supplementary Fig. 3, purple). Moreover, NgaP2 (Group 3-2) has two loops above and below the cleft from the second and eighth  $\beta$ -sheets of the ( $\beta/\alpha$ )<sub>8</sub> barrel domain (Fig. 3d and Supplementary Fig. 4, blue). The strict exo-type activity of Group 3-2 can be explained by these two loops that completely block the -2 subsite, forming a pocket-like architecture at the substrate-binding site and thus preventing substrate entry.

Similarly, BvGH123 in GH123 also has two additional  $\beta$ -hairpins and one extended loop from the second  $\beta$ -sheet of the ( $\beta/\alpha$ )<sub>8</sub> barrel domain, which blocks the -2 subsite side of the cleft, yielding a pocket-like conformation (Fig. 3e and Supplementary Fig. 5, magenta).

These findings revealed that substrate-binding sites with cleft structures exhibit endo-type activity, while those with pocket-like architectures show exo-type activity.

## Substrate recognition and catalytic mechanism of $\beta$ -GalNAc-acting enzymes

We further examined detailed substrate recognition mechanisms using structural analysis of complexes with GalNAc-thiazoline, an analog of the oxazolinium intermediate and a potent inhibitor of enzymes utilizing substrate-assisted catalysis (31, 32) (Fig. 4, Supplementary Fig. 6). With the exception of NgaCa, the structures of the NgaAt, NgaDssm, and NgaP2 complexes with GalNAc-thiazoline were successfully characterized (Supplementary Tables 2–4), and a subsequent docking model was engineered for NgaCa.

The docking model of NgaCa (Group 1) demonstrated that seven amino acid residues recognize GalNAc-thiazoline (Fig. 4a, green residues). In NgaAt (Group 2), apart from these seven residues, W329 establishes a hydrogen bond with GalNAc-thiazoline, and L555 contributes to substrate positioning (Fig. 4b, orange residues). Similarly, in NgaDssm (Group 3-1), the H479 residue forms a hydrogen bond, while residues W204 and L485 are involved in substrate positioning (Fig. 4c, purple residues). These results imply that Group 2 and Group 3-1 members possess a higher affinity for  $\beta$ -GalNAc than those of Group 1. Furthermore, complex binding modes were observed for Group 3-2 and GH123. For NgaP2 (Group 3-2), in addition to the seven residues discussed for NgaCa, three (D244, D297, and W298) and two (W196 and F489) residues are involved in hydrogen bond formation and substrate positioning, respectively (Fig. 4d, blue residues). Similarly, in BvGH123 (GH123), two (W253 and Q256) and three (W206, W306, and W482) residues are involved in these processes (Fig. 4e, magenta residues). Compared to the amino acid positions of the endo-type enzyme (green sticks), the exo-type enzymes have additional substrate-binding amino acid residues around  $\beta$ -GalNAc, indicating stronger recognition from the -2 subsite side. The increased number of residues could potentially reinforce substrate recognition in Group 3-2 and GH123 enzymes (Fig. 4f). Additional structural analyses (comparison of crystal structure with AlphaFold2 predicted structure of NgaCa, comparison of apo1-form and apo2-form of NgaCa, and comparison of GalNAc-thiazoline-bound form and GlcNAc-thiazoline-/apo-form of NgaAt, NgaDssm, and NgaP2) data are illustrated in Extended Data Fig. 7.

To complete the enzymatic analysis, we examined GalNAc- or GlcNAc-thiazoline as inhibitors of these enzymes. In the assays using GalNAc-thiazoline, Group 3-2 exo- $\beta$ -NGAs showed inhibitory activity at  $>1$  nM concentrations (Fig. 4g, blue and light blue). By contrast, no inhibition was observed in Group 1 endo- $\beta$ -NGA even at a concentration of 100  $\mu$ M (Fig. 4g, green). Interestingly, the activity of Group 3-1 endo/exo- $\beta$ -NGA was inhibited at concentrations 1000-fold higher than for exo- $\beta$ -NGA (Fig. 4g, purple). NgaAt was inhibited by 25 nM GalNAc-thiazoline and by 25  $\mu$ M GlcNAc-thiazoline (a 1000-fold higher concentration) (Fig. 4g, orange and yellow). These results corroborate the disparity in recognition abilities of GalNAc-thiazoline identified from the aforementioned structures.

### Analysis of catalytic residue mutants

Based on structural analysis, Asp and Glu of the DE motifs were recognized as the stabilizer of the 2-acetamido group and acid-base catalytic residue, respectively. The

GalNAc-thiazoline-bound structures and inhibition assays indicated that these enzymes perform substrate-assisted catalysis. Therefore, a point mutation analysis of the DE motifs was conducted (Fig. 4h). Mutants with Asp-to-Glu/Asn alterations exhibited little to no activity even at 10-fold to 100-fold higher enzyme concentrations. Glu-to-Asp/Gln mutants demonstrated reduced enzymatic activities, and the trends were similar to those observed for other enzymes performing substrate-assisted catalysis (9, 12).

### Elucidation of the catalytic mechanism using NMR

GH123 enzymes are exo- $\beta$ -NGAs, while enzymes from Groups 1 and 3-1 are endo- and endo/exo-type  $\beta$ -NGAs, respectively. To further confirm the substrate-assisted catalysis mechanism of these enzymes, we investigated the reaction products by NMR using Gal $\beta$ -1-3GalNAc $\beta$ -pNP as a substrate (Extended Data Figs. 8 and 9, Supplementary Figs. 7 and 8, Supplementary Tables 6 and 7) and monitored the stereochemistry of glycosidic bond hydrolysis (between the GalNAc and pNP moieties) by  $^1\text{H}$  NMR. The anomeric hydrogen signal of Gal $\beta$ -1-3GalNAc $\beta$ -pNP (between the GalNAc and pNP moieties) disappeared within 1 min for NgaDssm (Extended Data Fig. 8a) and 10 min for NgaCa (Extended Data Fig. 8b) and Gal $\beta$ -1-3GalNAc $\beta$  appeared, while Gal $\beta$ -1-3GalNAc $\alpha$  anomeric signals appeared after 10 and 30 min, respectively, due to mutarotation. These results indicate that NgaDssm and NgaCa are anomer-retaining enzymes, similar to GH123 enzymes.

### Discussion

We used deep-sea metagenomic sequences to discover a novel  $\beta$ -NGA, NgaDssm, which is the first  $\beta$ -NGA to possess dual endo/exo-type- $\beta$ -NGA activity and low sequence similarity to known exo- $\beta$ -NGAs. Prior studies have characterized three GH123 exo- $\beta$ -NGAs (NgaP, CpNga123, and BvGH123) from land soil (9) and human gut bacteria (12, 13); however, no genes from other ecological niches have been reported. The deep-sea environment (below 200 m depth) is characterized by total darkness, low temperatures, and high pressure, with occasional high temperatures owing to geological formations, such as hydrothermal vents. The deep-sea environments are completely distinct from the terrestrial environment and remain unexplored owing to limited accessibility to samples. Therefore, the deep-sea microbiome is an attractive potential bioresource for screening undiscovered enzymes (23).

We further discovered novel endo-, endo/exo-, and exo-type  $\beta$ -NGAs, as well as endo/exo- $\beta$ -HEXs, acting on  $\beta$ -GalNAc by analyzing deep-sea metagenomic sequences and public protein databases, and our comparative biochemical and structural analyses provide insights into the molecular evolution of  $\beta$ -GalNAc-targeting enzymes (Fig. 5). Despite their phylogenetic distance and different substrate specificities, these novel enzymes and GH123  $\beta$ -NGAs are likely homologous proteins based on the presence of conserved residues in their sequences (Fig. 1d), which are also positionally retained in their structures (Fig. 2c), and the observation that mutations in conserved residues were not tolerated and resulted in destabilization of protein structure and elimination of substrate recognition (Extended Data Figs. 6e, f). Structural comparisons between members in Groups 1–3 and GH123 enzymes further supported that the four families have diversified their substrate specificity (endo- $\beta$ -NGA for Group 1, endo/exo- $\beta$ -HEX for Group 2, endo/exo- $\beta$ -NGA, exo- $\beta$ -NGA for Group 3, and exo- $\beta$ -NGA for GH123) through the accumulation of point mutations and insertional sequences (Figs. 3 and 4). These data suggested a monophyletic evolutionary history of  $\beta$ -NGAs from the prototype enzymes in Group 1  $\beta$ -NGAs (Fig. 5).

In the CAZy database, several GH families are grouped into "clans" according to the conservation of catalytic residues, common reaction mechanisms, and structural similarities (11, 33). Consistent with these criteria, we suggest a new clan composed of the enzyme groups discovered in this study (Groups 1, 2, and 3) and GH123. Among the new enzymes identified herein, only NgaAt (Group 2) is a  $\beta$ -HEX that acts on  $\beta$ -GalNAc and  $\beta$ -GlcNAc. We speculate that a prototype enzyme of the new putative GH clan, including GH123 was a  $\beta$ -NGA specific for  $\beta$ -GalNAc, and the  $\beta$ -HEX enzymes of Group 2 are proposed to be divergent from the structurally simpler  $\beta$ -NGA. By contrast, GH20, mainly comprising  $\beta$ -HEX enzymes, belongs to clan GH-K, together with GH18 and GH85. All GH-K enzymes act on  $\beta$ -GlcNAc bonds (e.g., exo- and endo- $\beta$ -N-acetylglucosaminidase, chitinase, lacto-N-biosidase, and  $\beta$ -hyaluronidase). Thus, GH-K enzymes seem to originate from a structurally simple  $\beta$ -N-acetylglucosaminidase, and GH20  $\beta$ -HEX appears to belong to a family that diverged from the prototypical  $\beta$ -N-acetylglucosaminidase. Therefore, since  $\beta$ -HEX acts on  $\beta$ -GalNAc and  $\beta$ -GlcNAc, we suggest that  $\beta$ -HEX has two enzyme lineages with different evolutionary pathways.

In a previous study on  $\beta$ -NGAs from GH123, the function of DUF4091 has never been reported (12, 13). However, our domain-based exploration of  $\beta$ -NGA genes, based on GH123 members in the deep-sea metagenomic data, led to the identification of novel  $\beta$ -NGA types

containing DUF4091. Mutation experiments also highlighted that DUF4091 is indispensable for the stability of  $\beta$ -NGAs. The functional annotation of CDSs is typically conducted based on the similarity of the entire sequence to known protein sequences. A domain-based search (typically based on hidden Markov models) is an alternative method for identifying homologous sequence regions, even if the overall similarity is low. Examining CAZy families with consideration of co-occurring functionally unknown domains will likely lead to the subsequent discovery of additional novel enzymes.

To investigate the possible biological functions of the novel  $\beta$ -NGAs, we examined neighboring genes and protein-protein interactions using the STRING database (34) (Extended Data Fig. 10, Supplementary Data S4 and S5). The enzymes identified herein are potentially involved in the regulation of lipoprotein (NgaCa), the capsular membrane (NgaCp), and *O*-antigen (NgaCs), as well as in the degradation of glycans in the periplasm (NgaBf). Therefore,  $\beta$ -NGAs may exhibit more diverse functions than solely the degradation and utilization of glycans, which GH123  $\beta$ -NGA is thought to be responsible for and could have roles in various GalNAc-mediated biological processes.

Since the majority of prokaryotes in nature remain uncultured, it is crucial to explore the functional and genetic diversity of glycosidases in uncultured microorganisms through metagenomic analysis in order to gain a comprehensive understanding of glycan-mediated phenomena. The novel glycosidases hold potential for various industrial applications, including the utilization of these enzymes as biocatalysts for the production of functional oligosaccharides, glycan structure analysis, and disease diagnosis. Therefore, the discovery of novel glycosidases is crucial for both the basic and applied sciences. The comprehensive exploration and characterization of the diverse  $\beta$ -NGAs presented significantly enhance our understanding of their biological functions and their evolutionary history. These findings present a new approach to enunciating the evolutionary history of not only  $\beta$ -NGA but also many glycan-related enzymes. Further elucidation of the correlation between structures and diverse functions determined during evolution would also have significant implications for the structural basis of enzyme design for the engineering of new enzymes.



## Methods

### Sediment sampling and metagenomic sequencing

Abyssal sediment core was collected using a gravity corer on a remotely operated vehicle *ABISMO* (35), at station IOB (located 29.2746° N, 143.7673° E, at a depth of 5747 m below sea level) during a cruise of the ship R/V *Kairei* KR11-11 (December 2011) owned by the Japan Agency for Marine-Earth Science and Technology (JAMSTEC). The acquired sediment core was immediately subsampled onboard and stored at −80°C for molecular biology analyses (36).

Approximately 5 mL of frozen subsampled sediments from four sections (extracted from the depth ranges 0–8, 13–23, 53–63, and 113–123 cm below the seafloor) were used for metagenomic analysis. The sections were selected based on the geochemical profile and microbial community composition of the sediment core, as previously reported (36). Environmental DNA was extracted using the DNeasy PowerMax Soil Kit (QIAGEN, Hilden, Germany), per manufacturer protocol, with the following minor modification to increase DNA yield: cells were agitated twice for 10 min each after incubation at 65°C. Sequence libraries were prepared using the KAPA HyperPrep Kit (KAPA Biosystems) or the Ovation SP+ Ultralow Library System (NuGEN Technologies, San Carlos, CA, USA), as previously described (37). Library pools were mixed with Illumina PhiX control libraries and sequenced using the Illumina MiSeq or HiSeq platforms (Illumina, San Diego, CA, USA) at JAMSTEC or Macrogen (Seoul, South Korea).

### Bioinformatics

For raw metagenomic sequence data, both ends of the reads containing low-quality bases (Phread quality score <20) and adapter sequences were trimmed using TrimGalore (<https://github.com/FelixKrueger/TrimGalore>) with default settings. Sequencing reads derived from the PhiX genome were removed using Bowtie2 (38). Low complexity sequences or those shorter than 100 bp were discarded using PRINSEQ++ (39). The remaining high-quality paired-end reads of each sample were individually assembled *de novo* using metaSPAdes (40). Full-length coding sequences (CDSs) in the contigs were predicted using Prodigal (41) in the anonymous mode ('p meta' setting). Protein domain annotations of the CDSs were achieved through HMMER (42) against Pfam (version 35.0) (43) and dbCAN2 (version v10) (44) with a cutoff domain e-value of  $\leq 1E-3$ . All CDSs assigned to

GH123 and measuring > 400 bp in length were retrieved as  $\beta$ -NGA candidates and used for further analysis.

In addition to metagenomic CDSs, those with > 400 bp and architecture similar to known  $\beta$ -NGA genes (containing the DUF4091 [PF13320] domain in the C-terminal region) were retrieved from Pfam (version 35.0). A phylogenetic tree was constructed using MAFFT (45) with default settings and FastTree2 (46) with JC+CAT models. ClustalOmega (47) and the ColabFold software (48) were employed for sequence alignment and protein structure prediction, respectively.

### Construction of $\beta$ -NGA expression vectors

The selected  $\beta$ -NGA candidates were artificially synthesized by codon optimization for recombinant expression in *E. coli* using Strings DNA Fragments Synthesis Service (Thermo Fisher Scientific, Waltham, MA, USA) (Supplementary Data S3). The signal peptides in the amino acid sequences (Supplementary Data S1 and S3, highlighted with bold and underlined font) were predicted using SignalP 5.0 (49) and were removed from the N-terminal side of NgaDssm, NgaNp, NgaCs, NgaSc, and NgaCp. The gene encoding NgaP2 was directly cloned from the genomic DNA of *Paenibacillus* sp. TS12 via PCR amplification. These synthetic genes were designed with additional sequences at both terminals to facilitate amplification using a common primer set (Supplementary Data S6). The genes were cloned into the pET-47b(+) expression vector (Merck KGaA, Darmstadt, Germany) using the In-Fusion HD Cloning Kit (Takara Bio, Shiga, Japan). Mutagenesis was performed using the In-Fusion HD Cloning Kit. Primer sequences are listed in Supplementary Data S6.

### Expression and purification of recombinant $\beta$ -NGA

The expression vector and recombinant mutant plasmids were used to transform *E. coli* BL21 Star (DE3) cells. The cells were cultured in 50 mL of medium A (LB medium containing 50  $\mu$ g/mL of kanamycin), incubated at 37°C for 16 h with shaking, inoculated into 1–2 L of medium A, and incubated at 37°C for another 2–3 h with shaking. Protein expression was induced by addition of isopropyl  $\beta$ -D-1-thiogalactopyranoside (IPTG) to the culture at a final concentration of 0.1 mM. After additional culturing at 16°C for 16 h, cells were harvested by centrifugation (10,000 g for 10 min) and suspended in 50 mL of buffer A (20 mM HEPES-Na [pH 7.5], 150 mM NaCl, 5% [v/v] glycerol, 1 mM DTT, and 50 mM imidazole). Following sonication, cell debris was removed by centrifugation (13,000 g, 4 °C

for 20 min) and passed through a 0.45- $\mu$ m pore-sized GD/X syringe filter (Cytiva, Marlborough, MA, USA). The supernatant was subjected to chromatography using the AKTA Prime Chromatography System (Cytiva). The sample was loaded onto a 5 mL HisTrap HP column (Cytiva) at a flow rate of 2 mL/min. The column was then washed with buffer A. The His-tagged protein was eluted with Buffer B (20 mM HEPES-Na [pH 7.5], 150 mM NaCl, 5% [v/v] glycerol, 1 mM DTT, and 300 mM imidazole). The eluted fractions were pooled and dialyzed against Buffer C (20 mM HEPES-Na [pH 7.5], 150 mM NaCl, and 1 mM DTT). To cleave the His-Tag, the HRV3C protease was dialyzed at 4 °C for 16 h. The enzyme was further purified by ion-exchange [5-mL HiTrapQ column (Cytiva)] and size exclusion chromatography (SEC) [HiLoad 16/600 Superdex 200 pg column (Cytiva)]. The presence of the desired protein was confirmed by SDS-PAGE. The molecular weights of all  $\beta$ -NGAs estimated by SEC were consistent with the calculated molecular weights of their monomers.

## Enzyme assays

The activity of  $\beta$ -NGA candidates was determined using Assays I–III. In Assay I (*p*NP- $\beta$ -GalNAc as a substrate), the reaction mixture comprised 50 nmol of *p*NP- $\beta$ -GalNAc and an appropriate amount of the enzyme in 100  $\mu$ L of a 100 mM optimal buffer solution. Following a 30-min incubation, the reaction was arrested by adding 100  $\mu$ L of 1 M sodium carbonate, and the corresponding absorbance was measured at 405 nm. One unit of the enzyme was defined as the amount that catalyzed the release of 1  $\mu$ mol of *p*-nitrophenol per min from *p*NP- $\beta$ -GalNAc under experimental conditions. Values represent the mean of technical triplicate measurements. For Assay II (4MU- $\beta$ -GalNAc as a substrate), the reaction mixture was formulated using 10 nmol of 4MU- $\beta$ -GalNAc and an appropriate amount of enzyme in 100  $\mu$ L of a 100 mM buffer solution. Following incubation for 30 min, the fluorescence intensity was measured using a Synergy 2 multimode microplate reader (BioTek) at excitation and emission wavelengths of 360 and 460 nm, respectively. Values represent the mean of technical triplicate measurements. In Assay III (oligosaccharides as substrates), reaction mixtures containing 5 nmol of oligosaccharides and an appropriate amount of enzyme in 20  $\mu$ L of 100 mM buffer were incubated for 16 hours. The samples were boiled for 5 min to stop the reaction. The samples were dried, and the residues were dissolved in 10  $\mu$ L of a methanol:water (1:1, v/v) solution and applied to a TLC plate, which was then developed using a 1-butanol:acetic acid:water (2:1:1, v/v/v) solution. Oligosaccharides and

GalNAc were visualized using a diphenylamine-aniline-phosphate reagent (a mixture of 0.4 g diphenylamine, 0.4 mL aniline, 3 mL 85% phosphoric acid, and acetone [20 mL]). The optimal pH was determined using the GTA buffer [50 mM 3,3-dimethyl glutaric acid, 50 mM tris(hydroxymethyl)aminomethane, and 50 mM 2-amino-2-methyl-1,3-propanediol]. To avoid thermal degradation of *p*NP-substrates, the optimum temperatures of the enzymes were maintained between 0 and 75°C. The substrate specificity of the enzymes was examined using the following *p*NP-glycosides (50 nmol) or 4MU-glycosides (10 nmol): GalNAc-β- or α-*p*NP, Galβ1-3GalNAc-β- or α-*p*NP, GlcNAc-β- or α-*p*NP, galactose-β- or α-*p*NP, glucose-β- or α-*p*NP, arabinose-β- or α-*p*NP, mannose-β- or α-*p*NP, fucose-β- or α-*p*NP, xylose-β- or α-*p*NP, sulfate-*p*NP, GalNAc-β-4MU, GalNAc4S-β-4MU and GalNAc6S-β-4MU.

### Protein thermal shift assay

The protein thermal shift assay was conducted in the StepOnePlus Real-Time PCR System (Applied Biosystems) using Applied Biosystems Protein Thermal Shift Dye. The  $T_m$  of the proteins was calculated using the Protein Thermal Shift software v1.4 (Applied Biosystems).

### Crystallization, data collection, and structure determination

Purified β-NGA was concentrated to 10–20 mg/mL. This sample solution (0.5 μL) was mixed with an equivalent volume of the reservoir solution. Crystallization was performed by the sitting drop vapor diffusion technique at 20°C. Crystals were first formed in the crystallization screening trial and were reproduced by seeding the crystals into solutions prepared in an identical manner (Supplementary Data S7). For crystallization of the β-NGA-(GalNAc-thiazoline) complex, GalNAc-thiazoline was incorporated into the β-NGA protein solution to obtain a final concentration of 5 mM. GalNAc-thiazoline was prepared as previously described (31). A crystallization solution containing 20% (v/v) glycerol was used as a cryoprotectant for X-ray diffraction data collection. The X-ray diffraction experiments were performed at the BL32XU beamline of SPring-8. All diffraction data were collected using the automated data collection system ZOO (50). The obtained data were processed with XDS (51) using the automated data processing pipeline KAMO (52). For the X-ray diffraction data of NgaDssm complexed with GalNAc-thiazoline, automated structural analysis was performed using NABE (Matsuura et al., under review), and the structure was solved by molecular replacement using the AlphaFold2 model (48). PHENIX (53), COOT

(54), and REFMAC (55) were employed for structure refinement. Molecular images were displayed using PyMol (Schrödinger LLC, Palo Alto, CA, USA). The secondary structural elements in Supplementary Figs. 1–5 were determined using the ESPript software (56).

## Author Contributions

T. Sumida conceived and designed the study; performed molecular experiments, protein purifications, enzyme characterization, and protein crystallization; determined the protein structures; and wrote the manuscript. S.H. performed the bioinformatics analyses and wrote the manuscript. K.U. performed the molecular experiments and protein purifications. A.I. and K.T. performed NMR analysis. T. Sengoku performed structural prediction. K.A.S. synthesized inhibitors. S.D. and T.N. wrote the manuscript and supervised the project. S.F. determined the protein structures and wrote the manuscript. All authors reviewed the manuscript draft and approved the final manuscript.

## Acknowledgements

We would like to express sincere appreciation to the captain, crew, and all onboard scientists and technicians of the KR11-11 cruise. We are extremely grateful to the ROV *ABISMO* development and operation teams. Computations were partially performed on the NIG supercomputer at the ROIS National Institute of Genetics and the Data Analysis System and the Earth Simulator at JAMSTEC. This research was funded by the Research Support Project for Life Science and Drug Discovery (Basis for Supporting Innovative Drug Discovery and Life Science Research (BINDS)) from AMED under Grant Number JP22ama121001 (*support number 3118*). The synchrotron radiation experiments were performed at the BL32XU of SPring-8 with the approval of the Japan Synchrotron Radiation Research Institute (JASRI) (Proposal No. 2021A6700). We thank Kunio Hirata and BL32XU beamline staff for assisting with X-ray crystallographic data collection and analysis. We are grateful to Naohiro Matsugaki for the diffraction data analysis, Yoshitaka Moriwaki for the AlphaFold2 analysis, and Yukishige Ito for the NMR analysis. We thank Yasuhiro Shimane, Miho Hirai, and Fumie Kondo for their assistance in this study. This work was supported by the Japan Society for the Promotion of Science (Grant Numbers JP20K15444, JP22K05398) and the Mizutani Foundation for Glycoscience grant. K.A.S. appreciates the support of the Australian Research Council (FT100100291).

## Data availability

The metagenomic sequencing data were deposited in the DDBJ Sequence Read Archive under BioProject ID PRJDB15058 [<https://ddbj.nig.ac.jp/resource/bioproject/PRJDB15058>]. All coordinates are deposited in the PDB under accession numbers 8K2F, 8K2G, 8K2H, 8K2I, 8K2J, 8K2K, 8K2L, 8K2M, and 8K2N. All the requisite data for evaluating the conclusions are present in the article and/or the Materials and Methods. Additional information related to this study may be requested from the authors.

## References

1. Vinogradov, E., Sadovskaya, I., Grard, T. & Chapot-Chartier, M. P. Structural studies of the rhamnose-rich cell wall polysaccharide of *Lactobacillus casei* BL23. *Carbohydr. Res.* **435**, 156–161 (2016).
2. Hakomori, S. I. Structure and function of glycosphingolipids and sphingolipids: recollections and future trends. *Biochem. Biophys. Acta* **1780**, 325–46 (2008).
3. Hirata, T. *et al.* Identification of a Golgi GPI-*N*-acetylgalactosamine transferase with tandem transmembrane regions in the catalytic domain. *Nat. Commun.* **9**, 1–16 (2018).
4. Kelly, J. F. *et al.* Identification of a novel *N*-linked glycan on the archaellins and S-layer protein of the thermophilic methanogen, *methanothermococcus thermolithotrophicus*. *J. Biol. Chem.* **295**, 14618–14629 (2020).
5. Tailford, L. E., Crost, E. H., Kavanaugh, D. & Juge, N. Mucin glycan foraging in the human gut microbiome. *Front. Genet.* **5**, 1–18(2015).
6. Liu, B. *et al.* Structure and genetics of *Escherichia coli* O antigens. *FEMS Microbiol. Rev.* **44**, 655–683 (2020).
7. Badri, A. *et al.* Complete biosynthesis of a sulfated chondroitin in *Escherichia coli*. *Nat. Commun.* **12**, 1–10 (2021).
8. Varki, A *et al.* Essential of Glycobiology. 4<sup>th</sup> edition. Cold Spring Harbor (NY): Cold Spring Harbor Laboratory Press; (2022). doi:10.1101/9781621824213
9. Sumida, T., Fujimoto, K. & Ito, M. Molecular cloning and catalytic mechanism of a novel glycosphingolipid-degrading  $\beta$ -*N*-acetylgalactosaminidase from *Paenibacillus* sp. TS12. *J. Biol. Chem.* **286**, 14065–14072 (2011).
10. Henrissat, B. A classification of glycosyl hydrolases based on amino acid sequence similarities. *Biochem. J.* **280**, 309–316 (1991).

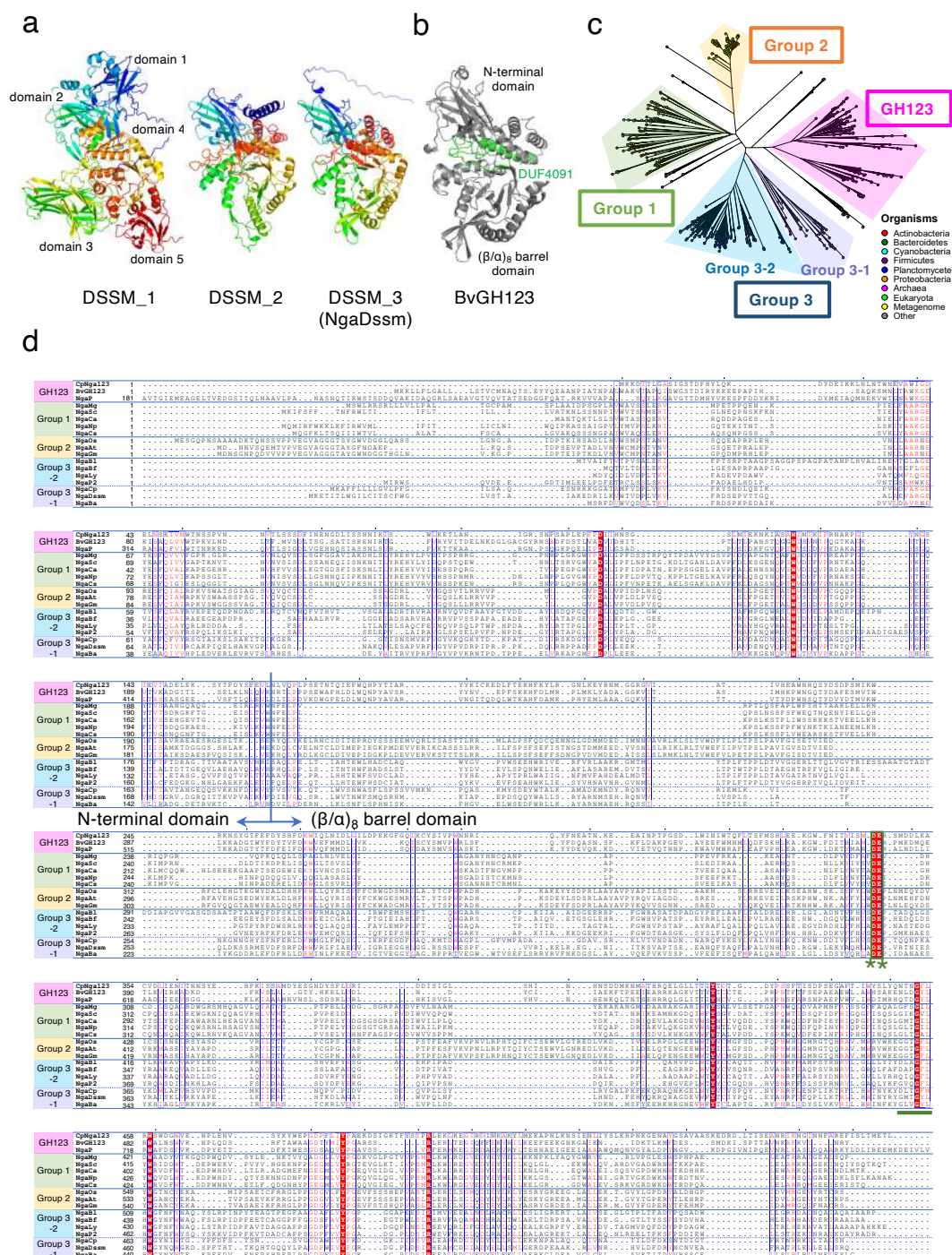


11. Henrissat, B., Bairoch, A. Updating the sequence-based classification of glycosyl hydrolases. *Biochem. J.* **316**, 695–696 (1996).
12. Noach, I. *et al.* The details of glycolipid glycan hydrolysis by the structural analysis of a family 123 glycoside hydrolase from *Clostridium perfringens*. *J. Mol. Chem.* **428**, 3253–3265 (2016).
13. Roth, C. *et al.* Structural and mechanistic insights into a *Bacteroides vulgatus* retaining N-acetyl- $\beta$ -galactosaminidase that uses neighbouring group participation. *Chem. Commun.* **52**, 11096–11099 (2016).
14. Drouillard, S., Armand, S., Davies, G. J., Vorgias, C. E. & Henrissat, B. *Serratia marcescens* chitobiase is a retaining glycosidase utilizing substrate acetamido group participation. *Biochem. J.* **328**, 945–949 (1997).
15. Lemieux, M. J. *et al.* Crystallographic Structure of Human  $\beta$ -Hexosaminidase A: Interpretation of Tay-Sachs Mutations and Loss of GM2 Ganglioside Hydrolysis. *J. Mol. Biol.* **359**, 913–929 (2006).
16. Sumida, T., Ishii, R., Yanagisawa, T., Yokoyama, S. & Ito, M. Molecular Cloning and Crystal Structural Analysis of a Novel  $\beta$ -N-Acetylhexosaminidase from *Paenibacillus* sp. TS12 Capable of Degrading Glycosphingolipids. *J. Mol. Biol.* **392**, 87–99 (2009).
17. Angelov, A. *et al.* A metagenome-derived thermostable  $\beta$ -glucanase with an unusual module architecture which defines the new glycoside hydrolase family GH148. *Sci. Rep.* **7**, 1–13 (2017).
18. Chuzel, L., Ganatra, M. B., Rapp, E., Henrissat, B. & Taron, C. H. Functional metagenomics identifies an exosialidase with an inverting catalytic mechanism that defines a new glycoside hydrolase family (GH156). *J. Biol. Chem.* **293**, 18138–18150 (2018).
19. Cheng, J. *et al.* Functional metagenomics reveals novel  $\beta$ -galactosidases not predictable from gene sequences. *PLoS One* **12**, 1–20 (2017).
20. Cabral, L. *et al.* Gut microbiome of the largest living rodent harbors unprecedented enzymatic systems to degrade plant polysaccharides. *Nat. Commun.* **13**, 1–16 (2022).
21. Strazzulli, A. *et al.* Discovery of hyperstable carbohydrate-active enzymes through metagenomics of extreme environments. *FEBS J.* **287**, 1116–1137 (2020).
22. Sogin, M. L. *et al.* Microbial diversity in the deep sea and the underexplored “rare biosphere”. *Proc. Natl. Acad. Sci. USA* **103**, 12115–20 (2006).
23. Tsudome, M. *et al.* An ultrasensitive nanofiber-based assay for enzymatic hydrolysis and deep-sea microbial degradation of cellulose. *iScience* **25**, 1–17 (2022).

24. Hiraoka, S. *et al.* Diverse DNA modification in marine prokaryotic and viral communities. *Nucleic Acids Res.* **50**, 1531–1550 (2022).
25. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
26. Strasser, R. *et al.* Enzymatic properties and subcellular localization of arabidopsis  $\beta$ -N-acetylhexosaminidases. *Plant Physiol.* **145**, 5–16 (2007).
27. Gutterigg, M. *et al.* Biosynthesis of truncated N-linked oligosaccharides results from non-orthologous hexosaminidase-mediated mechanisms in nematodes, plants, and insects. *J. Biol. Chem.* **282**, 27825–27840 (2007).
28. Kumar, A. B. *et al.* Tandem mass spectrometry has a larger analytical range than fluorescence assays of lysosomal enzymes: Application to newborn screening and diagnosis of mucopolysaccharidoses types II, IVA, and VI. *Clin. Chem.* **61**, 1363–1371 (2015).
29. Kumar, A. B. *et al.* Fluorimetric assays for N-acetylgalactosamine-6-sulfatase and arylsulfatase B based on the natural substrates for confirmation of mucopolysaccharidoses types IVA and VI. *Clin. Chim. Acta* **451**, 125–128 (2015).
30. Zhou, Z., Pan, J., Wang, F., Gu, J. D. & Li, M. Bathyarchaeota: Globally distributed metabolic generalists in anoxic environments. *FEMS Microbiol. Rev.* **42** 639–655 (2018).
31. Knapp, S & Myers D. S. Synthesis of  $\alpha$ -GalNAc Thioconjugates from an  $\alpha$ -GalNAc Mercaptan. *J. Org. Chem.* **67**, 2995–2999 (2002).
32. Sumida, T., Stubbs, K. A., Ito, M. & Yokoyama, S. Gaining insight into the inhibition of glycoside hydrolase family 20 exo- $\beta$ -N-acetylhexosaminidases using a structural approach. *Org. Biomol. Chem.* **10**, 2607–2612 (2012).
33. Davies, G.J. & Sinnott, M.L. Sorting the diverse: the sequence-based classifications of carbohydrate-active enzymes. *Biochem. J.* **30**, 26–32 (2008).
34. Szklarzyk, D. *et al.* The STRING database in 2021: customizable protein–protein networks, and functional characterization of user-uploaded gene/measurement sets. *Nucleic Acids Res.* **25**, D605–612 (2020).
35. Yoshida, H. *et al.* The ABISMO mud and water sampling ROV for surveys at 11,000 m depth. *Mar. Technol. Soc. J.* **43**, 87–96 (2009).
36. Hiraoka, S. *et al.* Microbial community and geochemical analyses of trans-trench sediments for understanding the roles of hadal environments. *ISME J.* **14**, 740–756 (2020).
37. Hirai, M. *et al.* Library construction from subnanogram DNA for pelagic sea water and deep-sea sediments. *Microbes Environ.* **32**, 336–343 (2017).

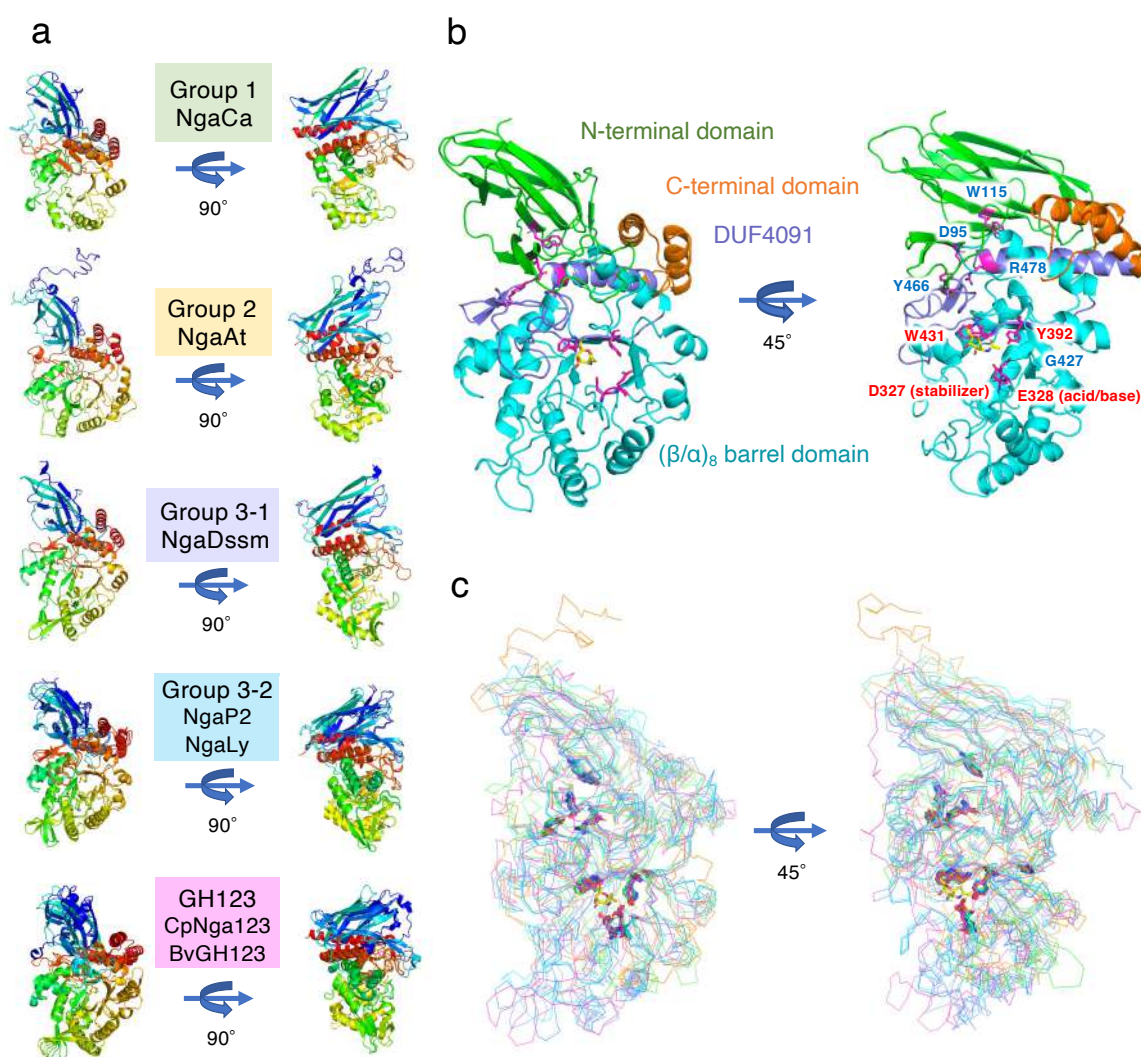
38. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
39. Cantu, V. A., Sadural, J. & Edwards, R. PRINSEQ++, a multi-threaded tool for fast and efficient quality control and preprocessing of sequencing datasets. *PeerJ Prepr.* **7**, e27553v1 (2019).
40. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
41. Hyatt, D. *et al.* Prodigal: Prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **11**, 119 (2010).
42. Mistry, J., Finn, R. D., Eddy, S. R., Bateman, A. & Punta, M. Challenges in homology search: HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Res.* **41**, e121 (2013).
43. Finn, R. D. *et al.* The Pfam protein families database: Towards a more sustainable future. *Nucleic Acids Res.* **44**, D279–D285 (2016).
44. Zhang, H. *et al.* dbCAN2: A meta server for automated carbohydrate-active enzyme annotation. *Nucleic Acids Res.* **46**, W95–W101 (2018).
45. Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
46. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2—Approximately maximum-likelihood trees for large alignments. *PLoS One* **5**, e9490 (2010).
47. Madeira, F. *et al.* Search and sequence analysis tools services from EMBL-EBI in 2022. *Nucleic Acids Res.* **50**, W276–W279 (2022).
48. Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nat. Methods* **19**, 679–682 (2022).
49. Almagro Armenteros, J. J. *et al.* SignalP 5.0 improves signal peptide predictions using deep neural networks. *Nat. Biotechnol.* **37**, 420–423 (2019).
50. Hirata, K. *et al.* Zoo: An automatic data-collection system for high-throughput structure analysis in protein microcrystallography. *Acta Crystallogr. D Struct. Biol.* **75**, 138–150 (2019).
51. Kabsch, W. XDS. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 125–132 (2010).
52. Yamashita, K., Hirata, K. & Yamamoto, M. KAMO: towards automated data processing for microcrystals. *Acta Crystallogr. D Struct. Biol.* **74**, 441–449 (2018).
53. Adams, P. D. *et al.* PHENIX: A comprehensive Python-based system for macromolecular structure solution. *Acta Crystallogr. D Biol. Crystallogr.* **66**, 213–221 (2010).

- 726 54. Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. Features and development of Coot. *Acta*  
727 *Crystallogr. D Biol. Crystallogr.* **66**, 486–501 (2010).
- 728 55. Murshudov, G. N. *et al.* REFMAC5 for the refinement of macromolecular crystal structures.  
729 *Acta Crystallogr. D Biol. Crystallogr.* **67**, 355–367 (2011).
- 730 56. Robert, X. and Gouet, P. Deciphering key features in protein structures with the new ENDscript  
731 server. *Nucleic Acids Res.* **42**, W320-4 (2014).



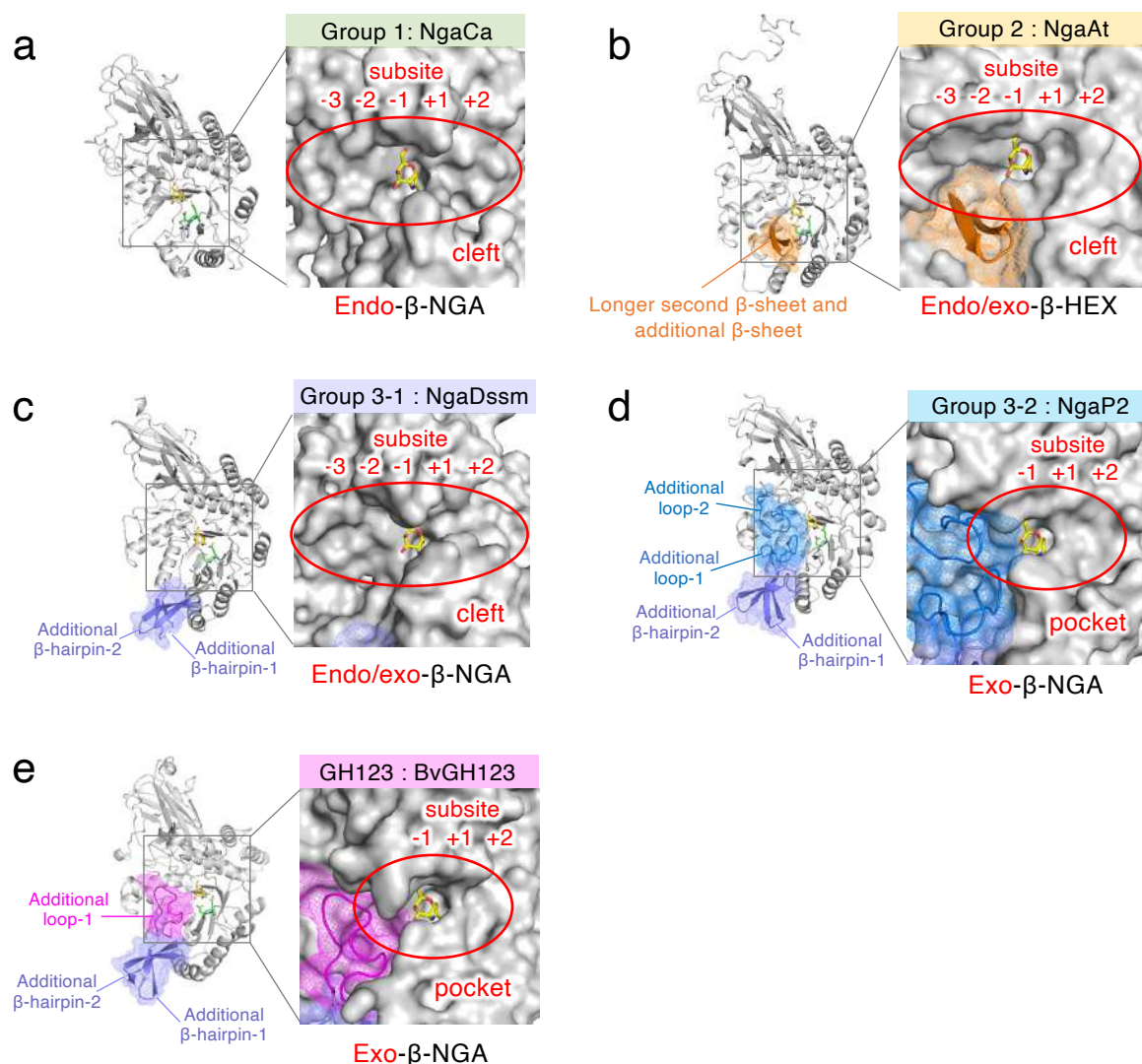
735 **a**, Overall structures of  $\beta$ -NGA candidates retrieved from deep-sea sediment metagenomes.  
736 The structures were predicted using AlphaFold2. **b**, The structure of BvGH123. DUF4091 is  
737 colored in green. **c**, Phylogenetic tree of  $\beta$ -NGA candidates retrieved from deep-sea sediment  
738 metagenomes and the Pfam database. **d**, Alignment of the  $\beta$ -NGA gene sequences. Residues  
739 conserved between all the analyzed proteins are shown (red background). The conserved DE  
740 residues (green asterisk, \*) are indicated by a green box. The DUF4091 region is underlined.





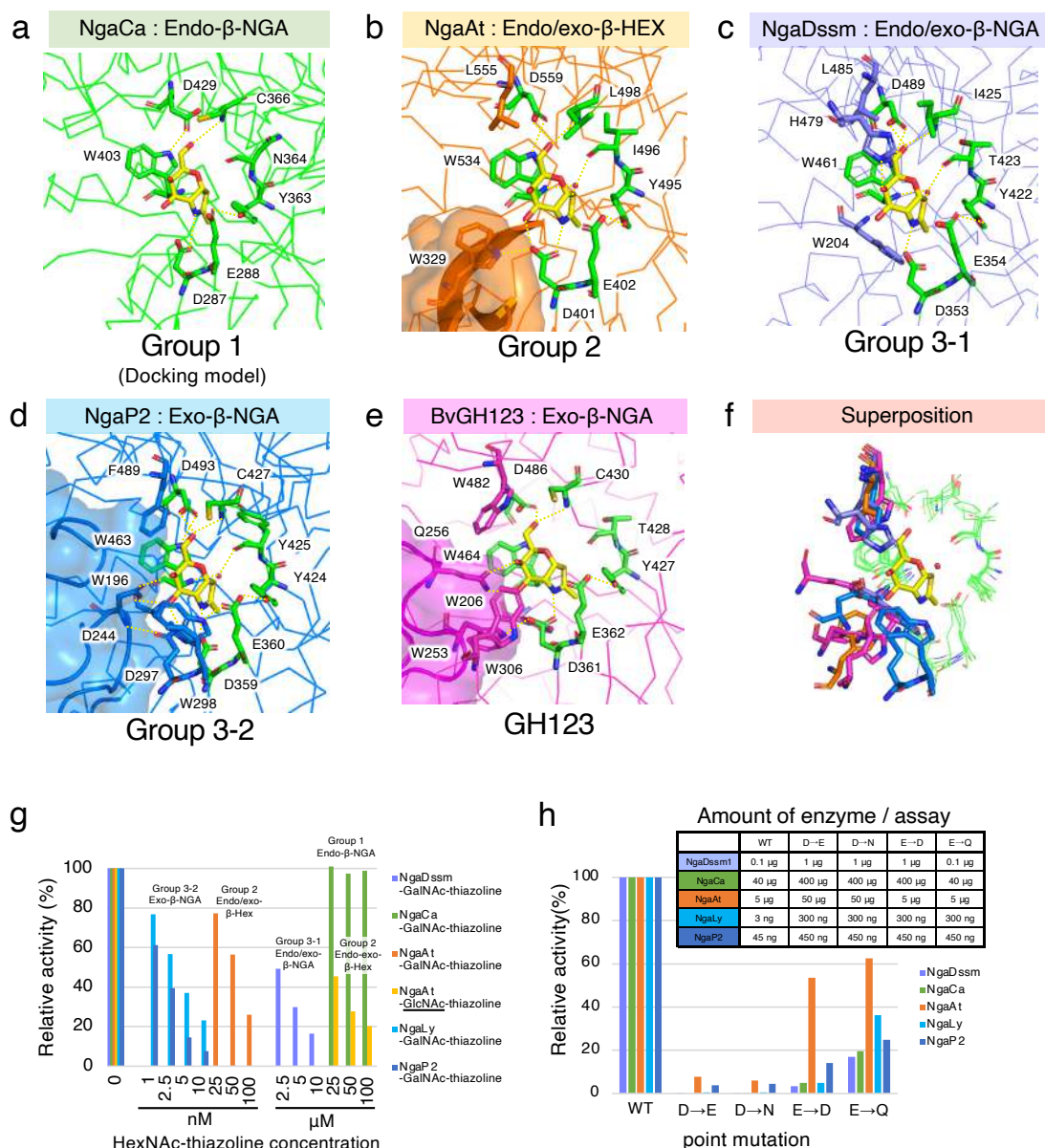
**Fig. 2. Overall structures and superposition of  $\beta$ -NGAs.**

**a**, Structures of  $\beta$ -NGAs. The DE motif is shown in stick format. **b**, The structures of NgaLy.  $\beta$ -Sandwich domain at the N-terminus,  $(\beta/\alpha)_8$ -barrel domain, DUF4091 and C-terminal domain are shown in green, cyan, purple, and orange, respectively. The conserved amino acids are shown as magenta sticks. Ligands are placed at the substrate-binding site and shown as yellow sticks. **c**, Superposition of the Group 1, Group 2, Group 3 and GH123 structures. Conserved amino acids are plotted on the structures as sticks.



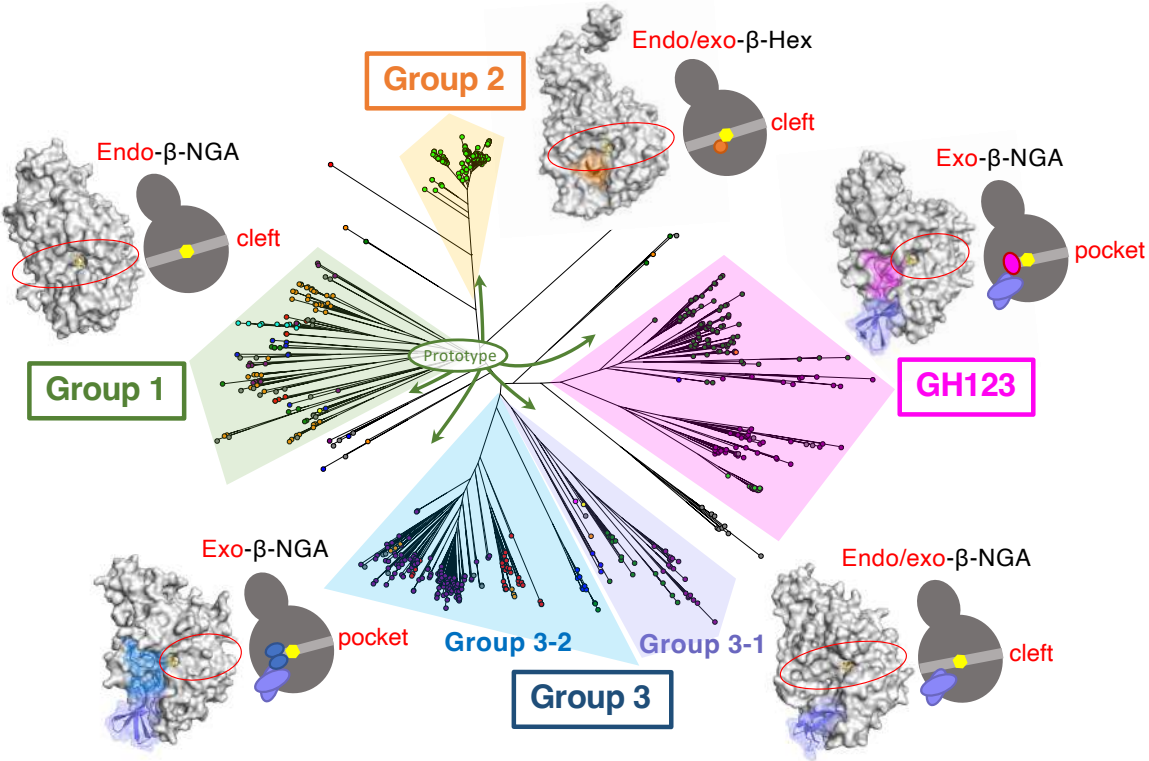
**Fig. 3. Substrate-binding area of β-NGAs.**

The overall structure of β-NGAs is shown as a cartoon (left) and the molecular surface of the substrate-binding site of β-NGA complexed with the ligand (right). The core structure (Group 1) is shown in gray (a), while additional regions characteristic of the other groups are shown in (b) orange (Group 2), (c) purple (Group 3-1), (d) blue (Group 3-2), and (e) magenta (GH123). The DE motif and ligand are indicated by green and yellow sticks, respectively.



**Fig. 4. Substrate recognition and catalytic mechanism.**

**a–e**, Active-site structure of each enzyme. **a**, Docking model of NgaCa with GalNAc-thiazoline. **b–e**, Crystal structure of each enzyme complexed with GalNAc-thiazoline. Hydrogen bonds are indicated by dotted lines. **f**, Superimposition of active sites in Groups 1, Group 2, Group 3 and GH123. **g**, Inhibition of  $\beta$ -NGA activity by GalNAc-thiazoline and GlcNAc-thiazoline. **h**, Point mutation analysis of the DE motif. The amounts of enzyme used in the assay are shown in the top-right table.



**Fig. 5. Phylogenetic tree and structural features of β-NGAs.**

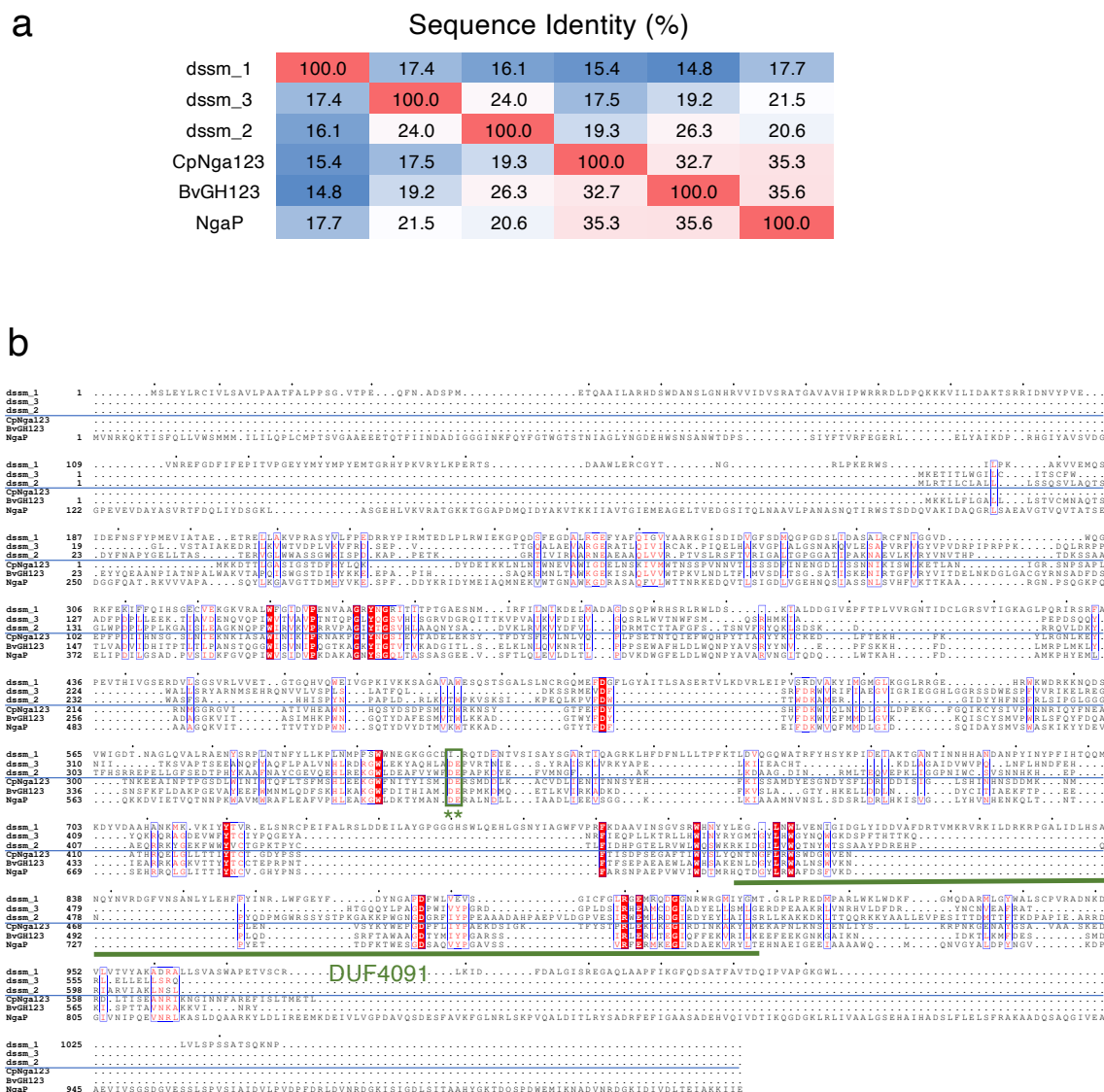
The overall structure of β-NGAs is illustrated as a surface model with an accompanying schematic representation. The ligand is indicated by a yellow stick. The Group 1 structure is depicted in gray as the basic structure. Additional regions characteristic of other groups are highlighted in orange (Group 2), purple (Group 3-1), blue (Group 3-2), and magenta (GH123).



**Table 1 Substrate specificities and general properties of  $\beta$ -NGAs**

	NgaDssm	NgaCa	NgaMg	NgaAt	NgaBa	NgaBf	NgaBl	NgaLy	NgaP2
	Group 3-1	Group 1	Group 1	Group 2	Group 3-1	Group 3-2	Group 3-2	Group 3-2	Group 3-2
Substrate specificity									
Substrate	Relative activity								
GalNAc- $\beta$ -pNP	100	-	-	100	100	100	100	100	100
Gal $\beta$ 1-3GalNAc- $\beta$ -pNP	70.3	100	-	-	7.3	-	-	-	-
Gal- $\beta$ -pNP	-	-	-	-	-	-	-	-	-
GlcNAc- $\beta$ -pNP	-	-	-	13.6	-	-	-	-	4.7
GalNAc- $\alpha$ -pNP	-	-	-	-	-	-	-	-	-
GalNAc- $\beta$ -4MU	54.0	-	-	100	100	100	100	100	100
GalNAc4S- $\beta$ -4MU	100	-	-	47.9	1.5	1.9	2.0	1.3	29.0
GalNAc6S- $\beta$ -4MU	-	-	-	-	-	-	-	-	-
General property									
Enzyme activity	Endo/exo- $\beta$ -NGA	Endo- $\beta$ -NGA	unknown	Endo/exo- $\beta$ -HEX	Endo/exo- $\beta$ -NGA	Exo- $\beta$ -NGA			
Optimal pH	5.0	5.0	-	5.0	6.5	6.5	5.0	7.0	5.5
Optimal buffer	Citrate	Citrate	-	Citrate	Citrate	Citrate	Citrate	HEPES	Citrate
Optimal temperature (°C)	45	25	-	40	70	35	45	40	40
$T_m$ (°C)	63.9	43.1	70.6	52.9	93.1	44.3	57.2	59.7	56/8
$K_m$ (mM)	3.8±0.42	10.5±2.7	-	12.4±4.6	3.9±0.54	0.98±0.069	2.8±0.44	0.98±0.052	0.36±0.025
$k_{cat}$ (s <sup>-1</sup> )	42.7±3.4	1.2±0.26	-	8.0±2.6	512.1±51.3	166.5±5.5	388.5±41.6	477.5±12.0	28.3±0.61
$k_{cat}/K_m$ (s <sup>-1</sup> mM <sup>-1</sup> )	11.2	0.11	-	0.64	130.3	170.1	139.7	489.2	78.5

- activity<0.5%. Values represent the means of technical triplicate experiments. General properties of each  $\beta$ -NGA are listed based on the results in Extended Data Fig. 3.



## Extended Data Fig. 1. Candidate $\beta$ -NGA sequences retrieved from deep-sea sediment metagenomes.

**a**, Sequence identities within  $\beta$ -NGA candidates retrieved from deep-sea sediment metagenomes and the known GH123 genes. **b**, Alignment of gene sequence. Residues conserved in all the proteins are shown on a red background. A conserved DE motif (green asterisk, \*) is indicated by a green box. The DUF4091 is underlined.



a

Sequence Identity (%)

GH123	CpNga123	100	33	35	16	19	17	17	18	15	16	16	12	14	13	14	16	18	17
	BvGH123	33	100	36	21	20	19	21	21	17	17	16	17	15	15	16	19	19	21
	NgaP	35	36	100	20	19	21	19	18	18	15	17	16	14	16	16	19	21	21
Group 1	NgaMg	16	21	20	100	39	41	40	40	24	23	24	18	18	18	17	20	19	26
	NgaSc	19	20	19	39	100	48	48	54	20	23	22	15	14	16	15	15	16	22
	NgaCa	17	19	21	41	48	100	55	53	22	23	23	18	19	18	21	18	21	24
	NgaNp	17	21	19	40	48	55	100	56	20	20	20	16	16	15	17	18	18	23
Group 2	NgaCs	18	21	18	40	54	53	56	100	22	22	21	17	17	16	18	19	19	25
	NgaOs	15	17	18	24	20	22	20	22	100	66	66	20	18	19	20	18	21	22
	NgaAt	16	17	15	23	23	23	20	22	66	100	74	19	17	20	21	17	20	23
Group 3 -2	NgaGm	16	16	17	24	22	23	20	21	66	74	100	20	18	20	19	18	20	23
	NgaBl	12	17	16	18	15	18	16	17	20	19	20	100	39	38	42	21	22	23
	NgaBf	14	15	14	18	14	19	16	17	18	17	18	39	100	40	39	20	25	23
Group 3 -1	NgaLy	13	15	16	18	16	18	15	16	19	20	20	38	40	100	40	22	24	23
	NgaP2	14	16	16	17	15	21	17	18	20	21	19	42	39	40	100	21	26	26
	NgaCp	16	19	19	20	15	18	18	19	18	17	18	21	20	22	21	100	36	29
Group 3 -1	NgaDsm	18	19	21	19	16	21	18	19	21	20	20	22	25	24	26	36	100	41
	NgaBa	17	21	21	26	22	24	23	25	22	23	23	23	23	26	29	41	100	

b

GH123	CpNga123	1	NTN	GLR	WSWDG	WVENP	...	LENVSY	...	KYWEF	GDFF	L	LY	FAEKDSIGKTF	YSTP	R	LEK	L	EGT	R	DINK	A	K	Y	L	M	
	BvGH123	1	NLD	GLR	WALNS	WVKNP	...	LQDSRF	...	TAWAA	GD	Y	M	Y	PGARSS	...	I	R	LE	R	L	T	E	G	T	F	
	NgaP	1	QTD	GLR	WAFDS	FVKDP	...	YETTD	...	KTWES	GD	S	A	Q	Y	P	CAVSS	...	I	R	P	E	R	M	K	E	
Group 1	NgaMg	1	G	IK	GLY	WRIDD	...	DEP	...	WEKVPVY	HGEKNFP	GE	C	M	LY	R	C	TEVGLKT	I	P	S	M	R	L	K	W	
	NgaNp	1	GLT	GLY	WQVD	WMT	...	KDP	...	WHDVLT	QTYSKNN	GD	N	F	PE	C	M	LY	P	C	Q	V	G	L	Q	G	
	NgaCs	1	GLT	GLY	WVRD	WMT	...	DDP	...	WNNVET	LQDGNHYP	GE	C	M	LY	P	C	Q	V	G	L	Q	G	I	E	G	
	NgaCa	1	GMT	GLY	WVRD	WMT	...	EDP	...	WHDVLT	LRADGMEFN	GE	C	M	LY	P	C	Q	V	G	L	Q	G	I	E	G	
Group 2	NgaOs	1	GGT	GLY	WGANC	YERAT	...	IPSAE	...	ICFRGL	PPDD	G	V	LY	P	C	E	V	F	S	S	S	H	E	P	A	
	NgaAt	1	GGT	GLY	WGANC	YERAT	...	VPSAE	...	VKFRRL	PPDD	G	V	LY	P	C	E	V	F	S	S	S	H	E	P	A	
	NgaGm	1	GGT	GLY	WGANC	YERAT	...	VASAE	...	IKFRHL	PPDD	G	V	LY	P	C	E	V	F	S	T	S	H	O	P	A	
Group 3 -2	NgaBl	1	DLA	GLY	WGFNF	YNAQY	...	SLRP	...	INPYTE	AGTPE	GF	FAA	GD	A	E	L	Y	P	C	P	...	GG	A	P	E	
	NgaBf	1	DAP	GLY	WGFNF	YNAQY	...	SLRP	...	IDPFE	ETCAG	GF	PF	GG	D	S	E	A	V	Y	P	C	P	...	EG	T	
	NgaLy	1	Q	Q	C	L	Y	W	G	F	N	F	Y	NAQ	L	S	T	R	P	D	P	F	A	V	T	D	
Group 3 -1	NgaP2	1	G	V	Q	C	L	Y	W	G	F	N	F	Y	NAQ	L	S	T	R	P	D	P	F	A	V	T	D
	NgaCp	1	DLT	GLY	WGFNF	YNAQY	...	SLRP	...	INPYTE	AGTPE	GF	FAA	GD	A	E	L	Y	P	C	P	...	GG	A	P	E	
	NgaDsm	1	GMT	GLY	WGFNF	YNAQY	...	SLRP	...	INPYTE	AGTPE	GF	FAA	GD	A	E	L	Y	P	C	P	...	GG	A	P	E	
	NgaBa	1	GLV	GLY	WGFNF	YNAQY	...	SLRP	...	INPYTE	AGTPE	GF	FAA	GD	A	E	L	Y	P	C	P	...	GG	A	P	E	

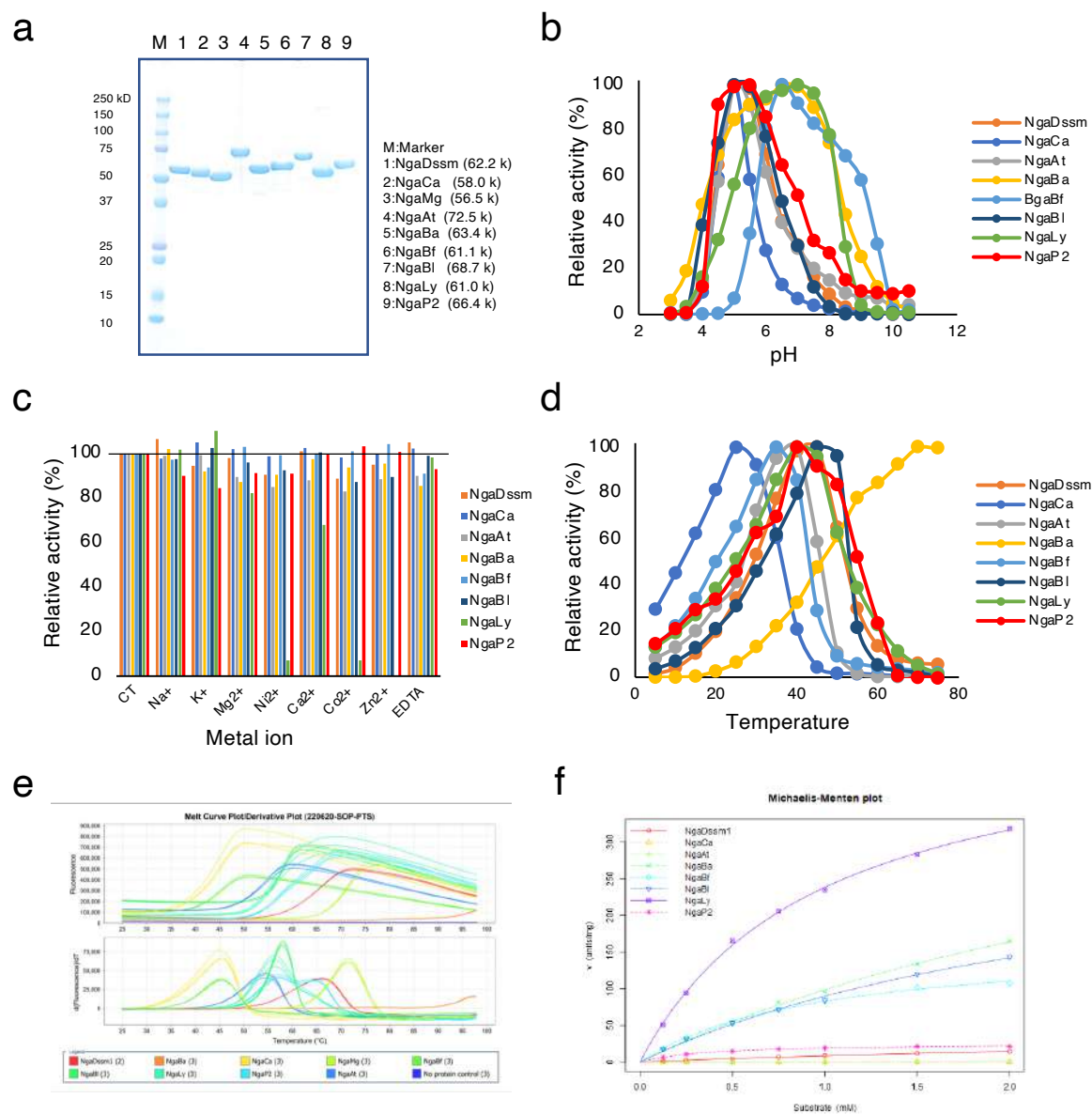
c

Sequence Identity (%) corresponding to DUF4091

GH123	CpNga123	100	43	46	29	26	23	29	27	30	30	27	25	22	25	28	31	25	34
	BvGH123	43	100	52	31	25	29	29	31	31	32	32	29	22	28	34	34	34	36
	NgaP	46	52	100	31	25	22	24	25	29	29	31	26	22	31	31	32	36	33
Group 1	NgaMg	29	31	31	100	51	59	59	55	30	34	31	21	21	24	26	30	28	40
	NgaSc	26	25	25	51	100	60	61	58	29	31	29	21	21	24	24	28	24	36
	NgaNp	23	29	22	59	60	100	69	67	29	33	30	22	23	25	25	33	30	40
	NgaCs	29	29	24	59	61	69	100	72	31	34	31	26	25	28	28	33	32	42
Group 2	NgaOs	27	31	25	55	58	67	72	100	32	32	32	25	28	31	34	33	38	38
	NgaAt	30	32	29	34	31	33	34	32	83	100	85	34	28	37	39	34	36	40
	NgaGm	27	32	31	31	29	30	31	32	79	85	100	32	28	35	34	33	34	37
Group 3 -2	NgaBl	25	29	26	21	21	22	26	25	31	34	32	100	51	49	51	35	32	35
	NgaBf	22	22	22	21	21	23	25	28	29	28	28	51	100	52	50	32	35	36
	NgaLy	25	28	31	24	24	25	28	31	34	37	35	49	52	100	67	32	38	38
Group 3 -1	NgaP2	28	34	31	26	24	25	28	34	36	39	34	51	50	67	100	35	42	46
	NgaCp	31	34	32	30	28	33	33	33	31	34	33	35	32	32	35	100	48	48
	NgaDsm	25	34	36	28	24	30	32	38	34	36	34	32	35	38	42	48	100	52
	NgaBa	34	36	33	40	36	40	42	38	37	40	37	35	36	38	46	48	52	100

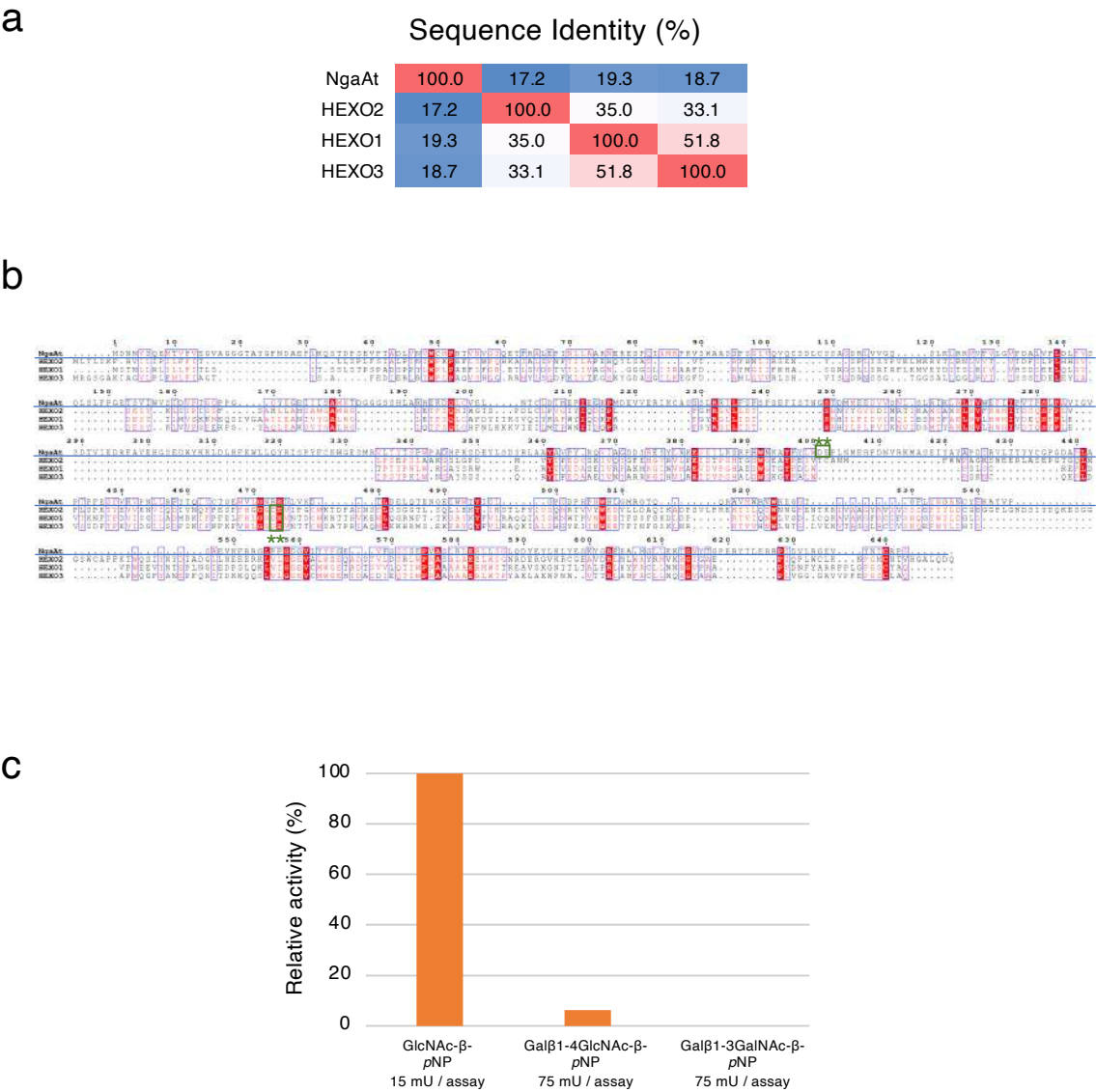
## Extended Data Fig. 2. Candidate $\beta$ -NGA sequences containing DUF4091.

a, Sequence identities within  $\beta$ -NGA genes. Names in red represent the enzymes whose activities were experimentally confirmed in this study. b, Alignment of the gene sequences corresponding to DUF4091. c, Sequence identities corresponding to DUF4091 with known GH123 and  $\beta$ -NGA candidates.



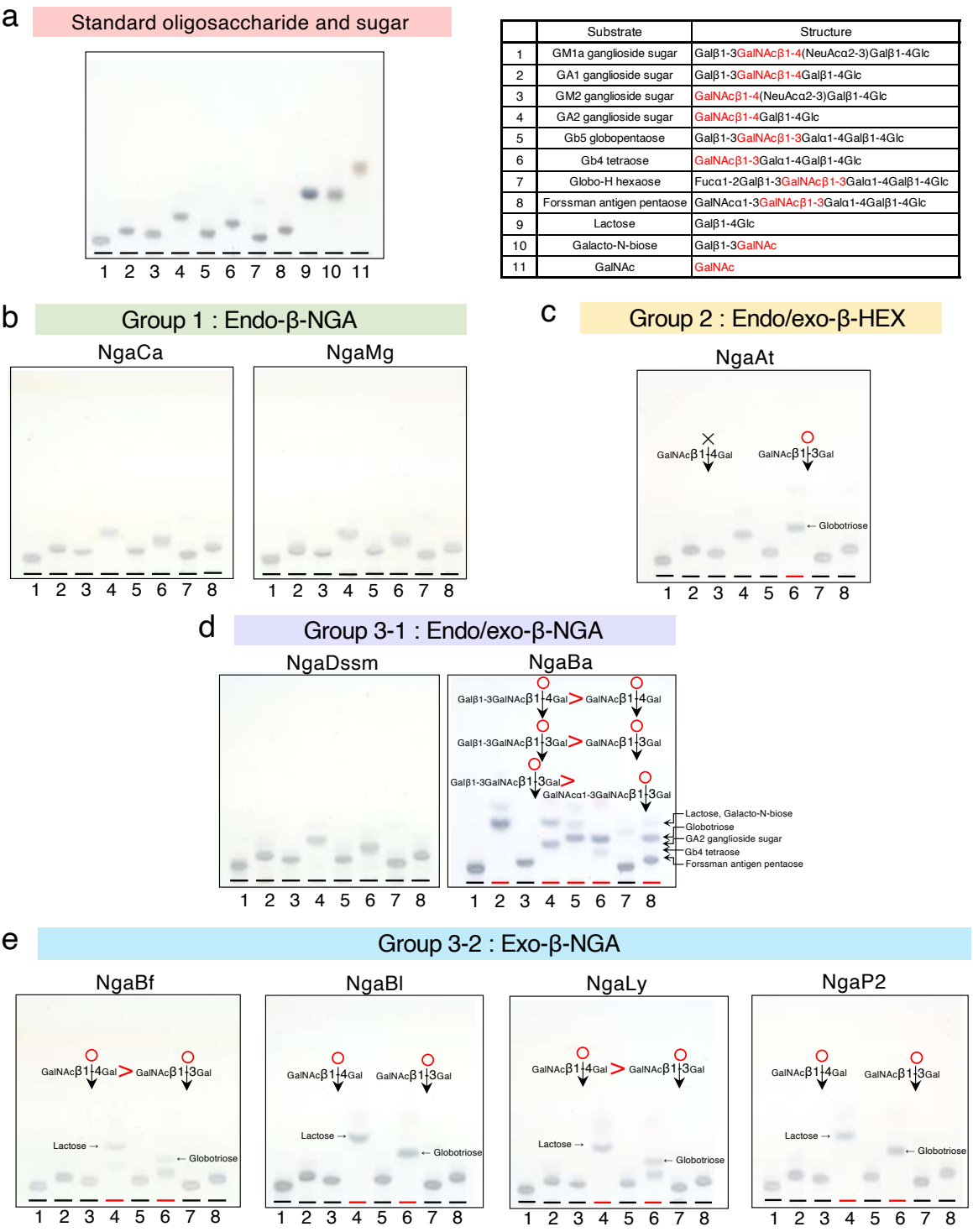
### Extended Data Fig. 3. General properties of recombinant β-NGAs.

**a**, SDS-PAGE analysis of recombinant β-NGAs, stained with Coomassie brilliant blue. Effect of **b**, pH and **c**, metal ions on enzymatic activity of recombinant β-NGAs at optimal temperature. **d**, Effect of temperature on enzymatic activity of recombinant β-NGAs. **e**, Protein thermal shift data of recombinant β-NGAs. **f**, s-v plots of recombinant β-NGAs. All values represent the mean of triplicate measurements.



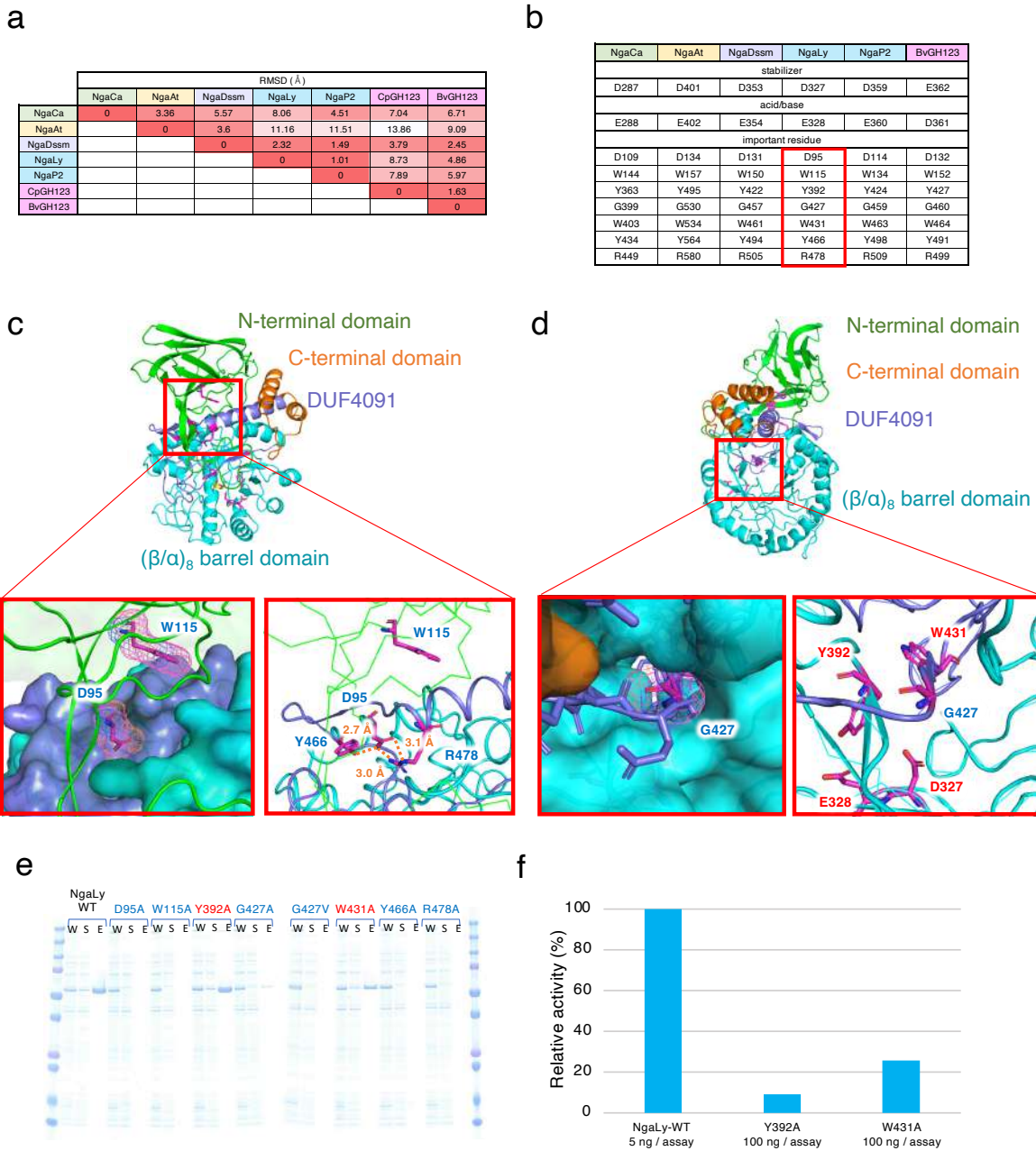
**Extended Data Fig. 4. Comparison between NgaAt and HEXO1–3.**

**a**, Sequence identities within the NgaAt and HEXO1–3 genes. **b**, Alignment of gene sequences. Residues conserved in all the proteins are displayed in a red background. The DE motifs (green asterisk, \*) of NgaAt and HEXO1–3 are indicated by green boxes. **c**, Specificity of NgaAt against different *pNP*-substrates. Values represent the mean of technical triplicate measurements.



**Extended Data Fig. 5. Hydrolysis of oligosaccharides by recombinant β-NGAs.**  
**a**, TLC analysis of the standard oligosaccharides and sugars (left), together with the structures

804 of the oligosaccharides used in the TLC assay (right). **b–e**, TLC demonstrates the hydrolysis  
805 of oligosaccharides of GM1a, asialo GM1 (GA1), GM2, asialo GM2 (GA2), Gb5, Gb4,  
806 Globo H, and Forssmann antigens. Oligosaccharides in each lane were arranged in the same  
807 order as in the standard TLC.

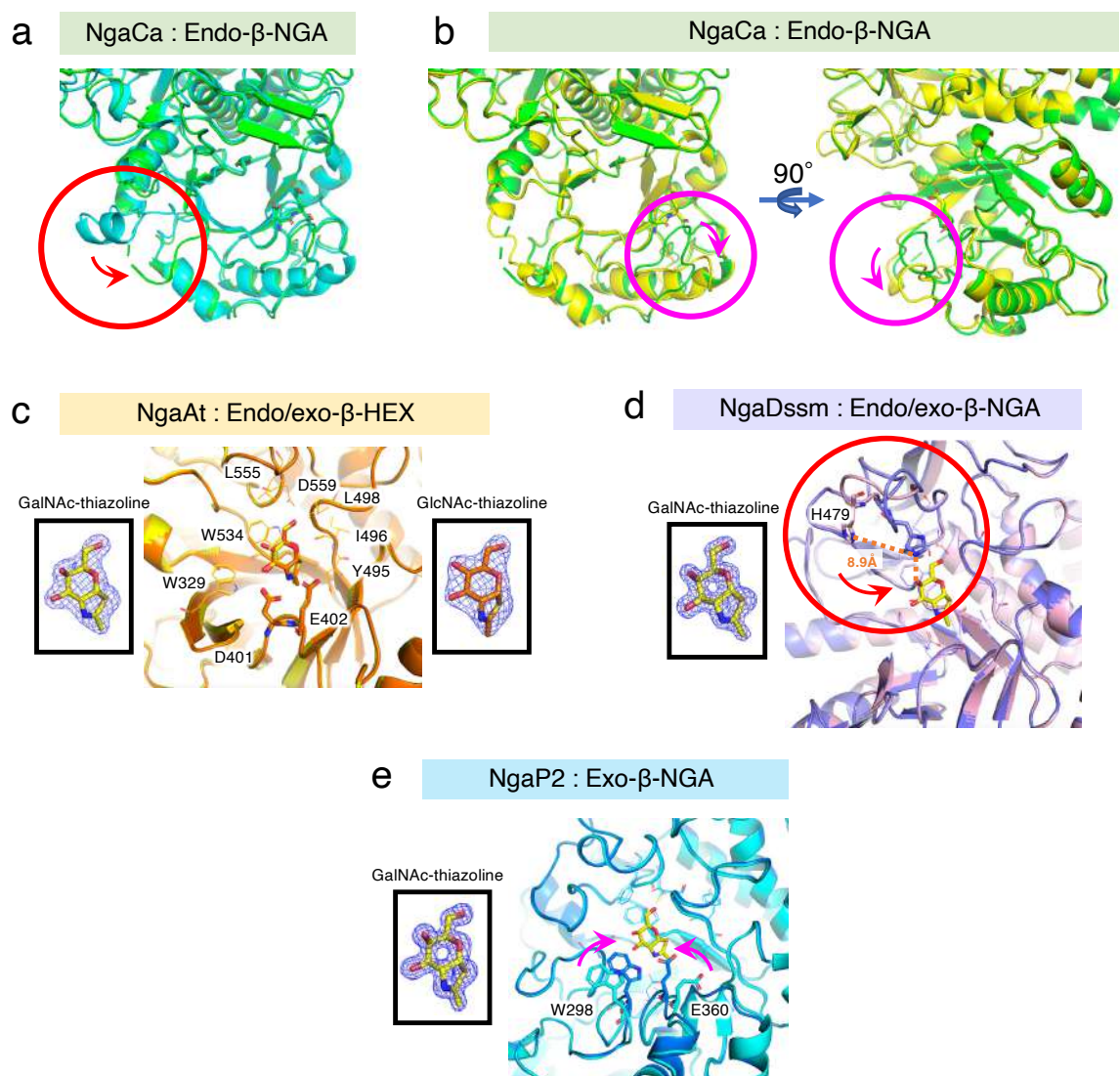


**Extended Data Fig. 6. Comparison of overall structure and effect of the mutation on amino acids conserved throughout all β-NGA groups.**

**a**, Structural similarity between β-NGA groups compared using the root mean square distance (RMSD). **b**, Amino acids conserved between β-NGA groups. Point mutation experiments were performed on the amino acids of NgaLy as specified in the red box. **c**, **d**, Detailed views of interactions between conserved residues (**c**, D95, W115, Y466, R478, **d**, G427). The β-



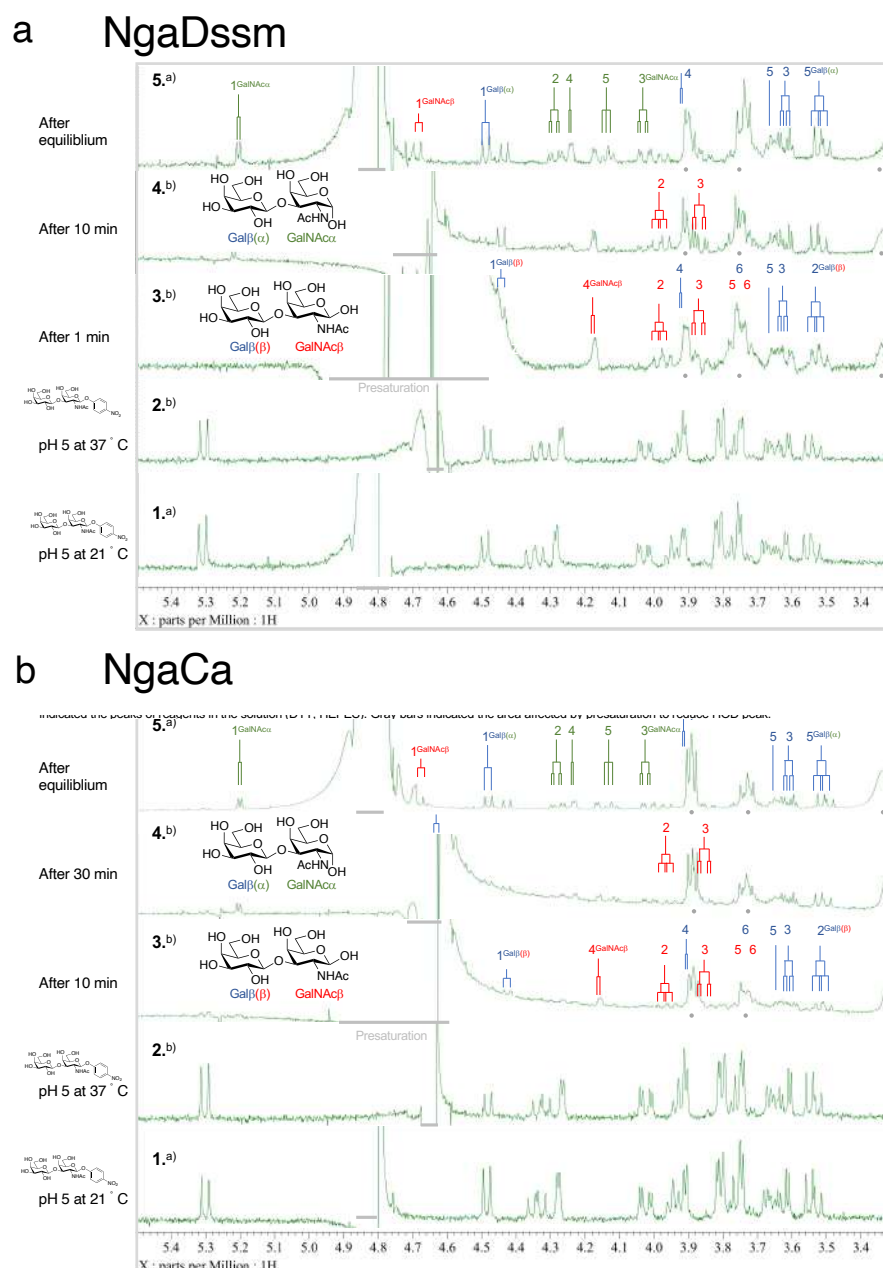
815 sandwich domain at the N-terminus, the ( $\beta/\alpha$ )<sub>8</sub>-barrel domain, DUF4091, and the C-terminal  
816 domain are depicted in green, cyan, purple, and orange, respectively. Magenta sticks  
817 represent the conserved amino acids. The ligands are indicated by yellow sticks. **e**, SDS-  
818 PAGE analysis of point mutants (W: whole-cell lysate; S: supernatant protein after sonication  
819 and centrifugation; E: eluted protein after affinity purification). The amino acids that likely  
820 contribute to structural stability and those involved in substrate recognition are demarcated  
821 with blue and red colors, respectively. **f**, Relative activity of the NgaLy Y392A and W431A  
822 mutants compared to that of the wild-type (WT) enzyme. Values represent the mean of  
823 technical triplicate measurements.



# **Extended Data Fig. 7. The structures of NgaCa, NgaAt, NgaDssm, and NgaP2.**

**a**, Comparison between the apo1 crystal structure (green) and AlphaFold2 predicted structure (cyan) of NgaCa. The position of the loop behind the second  $\beta$ -sheet in the  $(\beta/\alpha)_8$ -barrel domain is demarcated with a red circle. The DE motif is indicated by a stick. **b**, Distinction between apo 1 (green) and apo 2 (yellow) in NgaCa. The apo 2 structure is derived from crystals to which 5 mM Gal $\beta$ 1-3GalNAc was incorporated during crystallization, but Gal $\beta$ 1-3GalNAc is not visible. However, the loops surrounding the active site are different from those surrounding apo 1, and the loops move as the active site cleft expands (magenta circle). These two distinct states were designated as closed (apo 1) and open (apo 2), respectively.

834 The DE motif is indicated by a stick. **c**, Superimposition of the GalNAc- and GlcNAc-  
835 thiazoline bond forms of NgaAt (yellow and orange, respectively). The positions of the amino  
836 acids involved in substrate recognition are identical. GalNAc- and GlcNAc-thiazoline are  
837 represented by yellow and orange sticks, respectively. Polder maps of GalNAc- and GlcNAc-  
838 thiazoline ( $4\sigma$ ) are illustrated as blue meshes. **d**, Superimposition of the apo (pink) and the  
839 GalNAc-thiazoline-bound forms (purple) of NgaDssm. In the apo form, His479 is distant  
840 from the active site. In the GalNAc-thiazoline bound, His 479 was located within hydrogen-  
841 bonding distance from the 4-OH of GalNAc-thiazoline. These two distinct forms are  
842 designated as closed and open states, respectively. GalNAc-thiazoline is depicted by a yellow  
843 stick. A polder map of GalNAc-thiazoline ( $4\sigma$ ) is displayed as a blue mesh. **e**,  
844 Superimposition of the apo form (cyan) of NgaP2 onto its GalNAc-thiazoline-bound form  
845 (blue). The residues that shifted upon GalNAc-thiazoline binding, Trp298 and Glu360, are  
846 depicted as sticks. The displacements of these residues are indicated by arrows. GalNAc-  
847 thiazoline is represented by a yellow stick. A polder map of GalNAc-thiazoline ( $4\sigma$ ) is  
848 illustrated as a blue mesh.

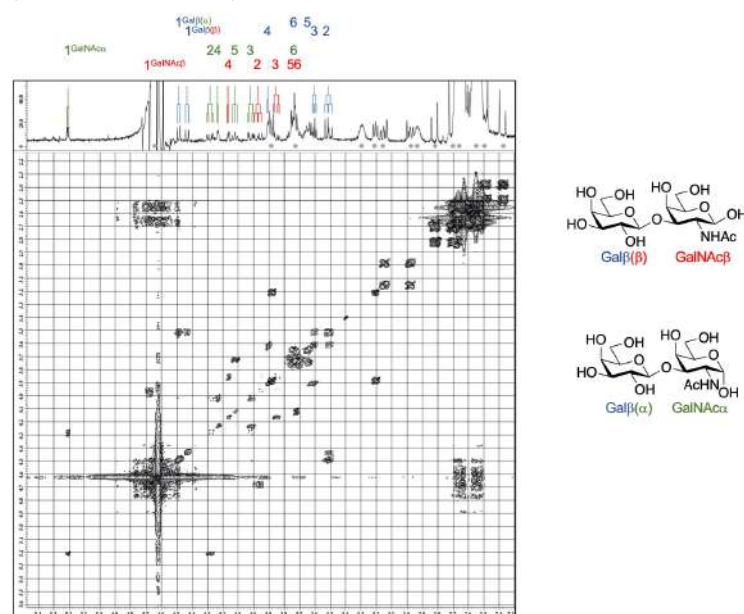


# **Extended Data Fig. 8. NgaDssm and NgaCa are anomer-retaining enzymes.**

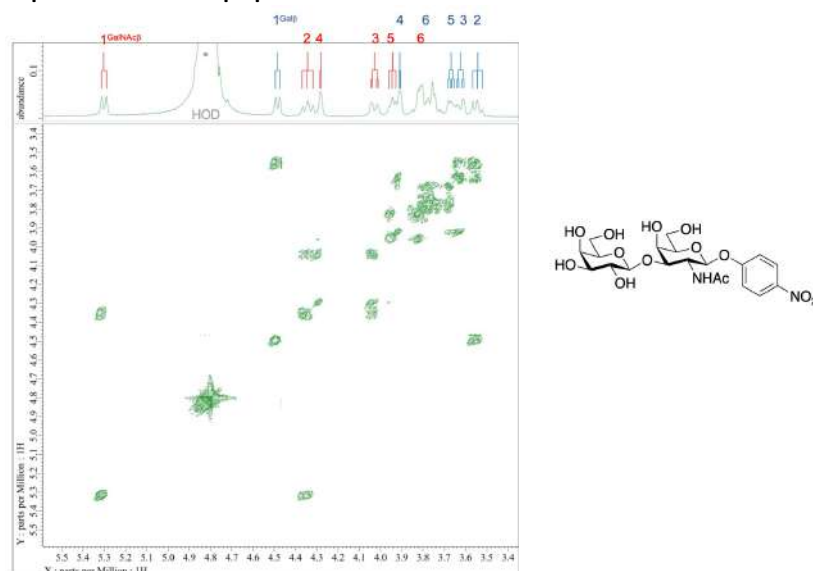
**a**,  $^1\text{H}$  NMR spectrum monitoring the activity of NgaDssm toward Gal $\beta$ 1-3GalNAc- $\beta$ -pNP in  $\text{D}_2\text{O}/\text{H}_2\text{O}$  (32:1, 100 mM citrate buffer, pH(D) 5.0) at 37°C. 1: Substrate at 21°C; 2: Substrate at 37°C at pH(D) 5.0 before addition of enzyme solution; 3: Reaction after 1 min; 4: Reaction after 10 min; 5: Reaction after attaining equilibrium. The enzyme was dissolved in 20 mM HEPES-Na (pH 7.5), 150 mM NaCl, and 1 mM DTT in  $\text{H}_2\text{O}$  and premixed with  $\text{D}_2\text{O}$

856 (D<sub>2</sub>O/H<sub>2</sub>O = 6:1) before treatment with the substrate. **b**, <sup>1</sup>H NMR spectrum monitoring the  
 857 activity of NgaCa toward Galβ1-3GalNAc-β-*p*NP in D<sub>2</sub>O/H<sub>2</sub>O (11:1, 100 mM citrate buffer,  
 858 pH(D) 5.0) at 37°C. 1: Substrate at 21°C; 2: Substrate at 37°C D<sub>2</sub>O/H<sub>2</sub>O (11:1, 100 mM citrate  
 859 buffer, pH(D) 5.0); 3: Reaction after 10 min; 4: Reaction after 30 min; 5: Reaction after  
 860 attaining equilibrium. The enzyme was dissolved in 20 mM HEPES-Na (pH 7.5), 150 mM  
 861 NaCl, and 1 mM DTT in H<sub>2</sub>O (32.5 mg/mL) and premixed with D<sub>2</sub>O (D<sub>2</sub>O/H<sub>2</sub>O = 2:1) before  
 862 treatment with the substrate. a) Measured at 21°C with presaturation at 4.80 ppm. b)  
 863 Measured at 37°C with presaturation at 4.63 ppm. The gray dots symbolize the peaks of the  
 864 reagents in solution (DTT and HEPES). The gray bars indicate the areas affected by  
 865 presaturation to reduce the HOD peak.

a Gal $\beta$ 1-3GalNAc $\alpha$ / $\beta$



b Gal $\beta$ 1-3GalNAc- $\beta$ -pNP

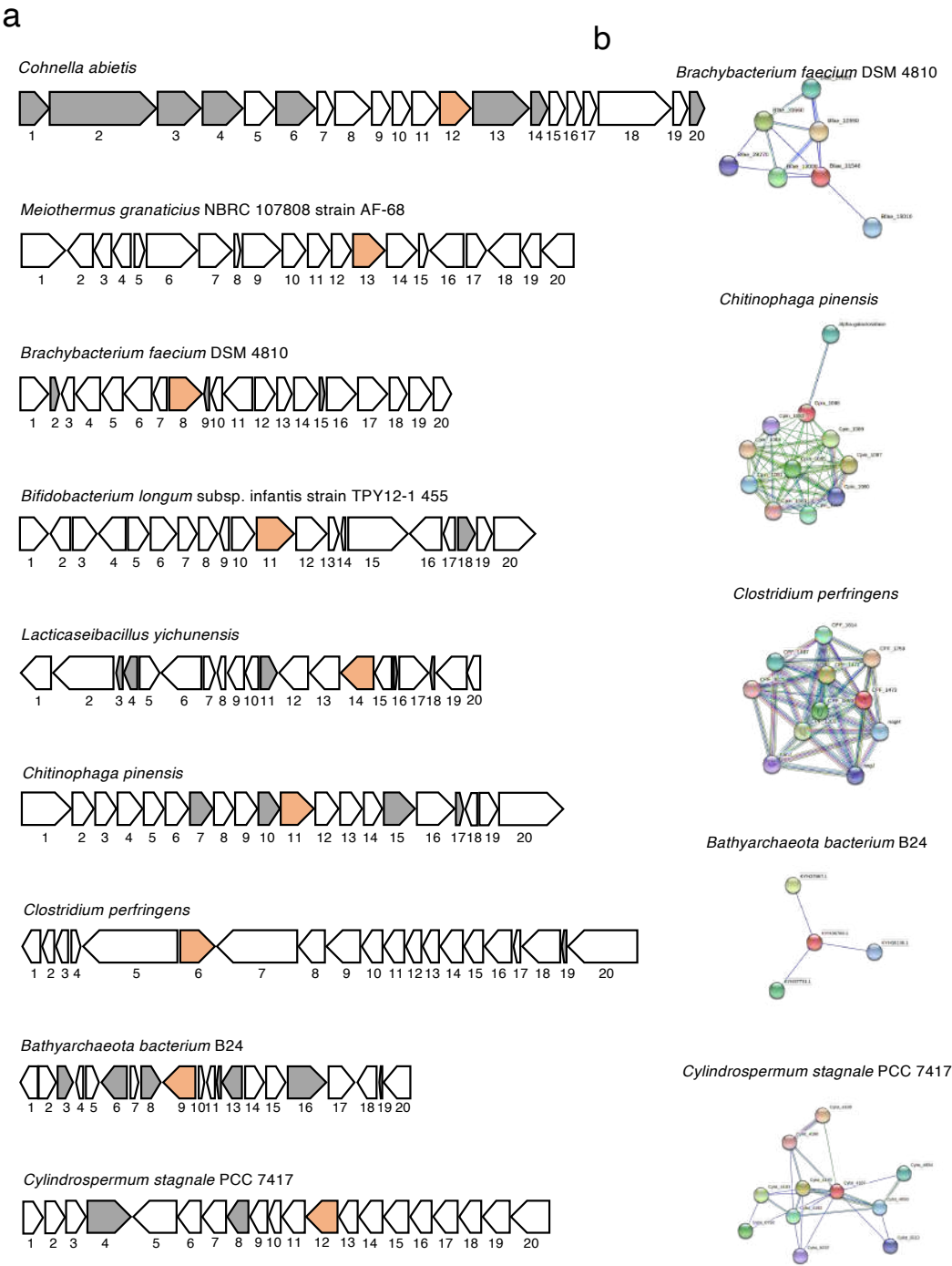


**Extended Data Fig. 9 H-H COSY and  $^1\text{H}$  NMR spectra of Gal $\beta$ 1-3GalNAc $\alpha$ / $\beta$  and Gal $\beta$ 1-3GalNAc- $\beta$ -pNP.**

**a**, 2D H-H COSY and 1D  $^1\text{H}$  NMR spectra of the reaction mixtures at equilibrium after hydrolysis. The mixture comprises Gal $\beta$ 1-3GalNAc $\alpha$ / $\beta$  (~1.5:1). Assignments were carried out by 1D  $^1\text{H}$  NMR and 2D H-H COSY as the chemical shift of  $^1\text{H}$  at 1 of  $\beta$ -GalNAc (4.58



872 ppm) is overlapped under the large HOD peak (4.4–4.8).  $\alpha$  and  $\beta$  in parentheses indicate the  
 873 stereochemistry of the reducing GalNAc residue of the corresponding disaccharides. **b**, H-H  
 874 COSY and 1D  $^1\text{H}$  NMR spectra of Gal $\beta$ 1-3GalNAc- $\beta$ -*p*NP in D<sub>2</sub>O. Assignments were  
 875 performed using 1D  $^1\text{H}$  and  $^{13}\text{C}$ , 2D H-H COSY, TOCSY, HMQC, HMQC-TOCSY, and  
 876 HMBC.



**Extended Data Fig. 10. Neighborhood genes and potential protein interaction networks of  $\beta$ -NGAs.**

**a**, Neighborhood genes around the  $\beta$ -NGA. Functional annotations were obtained from the

881 NCBI GenBank database, as listed in Supplementary Data S4. The  $\beta$ -NGAs and hypothetical  
882 proteins are depicted in orange and gray, respectively. **b**, Interaction networks of  $\beta$ -NGA  
883 retrieved from the STRING database (v11.5). Predicted functional partner proteins are listed  
884 in Supplementary Data S5.