

¹ Diversity of sugar-diphospholipid-utilizing glycosyltransferase families

² Ida K.S. Meitil¹, Garry P. Gippert¹, Kristian Barrett¹, Cameron J. Hunt¹, Bernard Henrissat^{1,2,3*}

³ December 10, 2023

⁴ **Abstract**

⁵ Peptidoglycan polymerases, enterobacterial common antigen polymerases, O-antigen ligases, and other bacte-
⁶ rial polysaccharide polymerases (BP-Pol) are glycosyltransferases (GT) that build bacterial surface polysac-
⁷ charides. These integral membrane enzymes share the particularity of using diphospholipid-activated sugars
⁸ and were previously missing in the carbohydrate-active enzymes database (CAZy; www.cazy.org). While the
⁹ first three classes formed well-defined families of similar proteins, the sequences of BP-Pols were so diverse
¹⁰ that a single family could not be built. To address this, we developed a new clustering method where a
¹¹ sequence similarity network was used to define small groups of alignable sequences, hidden Markov models
¹² (HMMs) were built for each group, and the resulting HMMs were aligned to form new families. Overall, we
¹³ have defined 17 new GT families including 14 of BP-Pols. We find that the reaction stereochemistry appears
¹⁴ to be conserved in each of the defined BP-Pol families, and that the BP-Pols within the families transfer
¹⁵ similar sugars even across Gram-negative and Gram-positive bacteria. Comparison of the new GT families
¹⁶ reveals three clans of distantly related families, which also conserve the reaction stereochemistry.
¹⁷

1 Introduction

¹⁸ Carbohydrate polymers (glycans) and glyco-conjugates are the most abundant biomolecules on Earth and
¹⁹ adopt a wide range of functions including energy storage, structure, signaling, and mediators of host-pathogen
²⁰ interactions [1]. Due to the stereochemical diversity of monosaccharides and the many possible linkages they
²¹ can engage into, glycans display an enormous structural diversity [2, 3]. Yet, our knowledge on their assembly
²² is far from complete, especially in comparison to the enzymes catalyzing their breakdown.

²³ The transfer of sugar moieties to acceptor molecules such as proteins, lipids or other sugars, is catalyzed
²⁴ by enzymes called glycosyltransferases or GTs [4]. Campbell and colleagues (1997) proposed a sequence-based
²⁵ classification of GTs into 26 families. The number of sequence-based families has since continued to grow based
²⁶ on the necessary presence of at least one experimentally characterized founding member to define a family, and
²⁷ is presented in the carbohydrate-active enzymes database (CAZy; www.cazy.org) [5]. An advantage of the
²⁸ sequence-based classification is that it readily enables genome mining for the presence of new family members.
²⁹ Today there are 116 GT families in the CAZy database and in contrast to the EC numbers [6], the sequence-
³⁰ based classification implicitly incorporates the structural features of GTs including the conservation of the
³¹ catalytic residues.

³² It was recognized very early that sequence-based GT families group together enzymes that can utilize
³³ different sugar donors and/or acceptors, illustrating how GTs can evolve to adopt novel substrates and form
³⁴ novel products [7, 8]. Mechanistically, glycosyltransferases can be either retaining or inverting, based on the
³⁵ relative stereochemistry of the anomeric carbon of the sugar donor and of the formed glycosidic bond [4].
³⁶ With almost no exceptions, this feature is conserved in previously defined sequence-based families, providing
³⁷ predictive power to this classification, as the orientation of the glycosidic bond can be predicted even if the
³⁸ precise transferred carbohydrate is not known.

³⁹ The large majority of the 116 GT CAZy families use donors activated by nucleotide diphosphates. Eleven
⁴⁰ families utilize nucleotide monophospho-sugars (sialyl and KDO transferases), while 12 families utilize lipid
⁴¹ monophospho-sugars. Until now, only one family in the CAZy database utilizes sugar-diphospholipid donors:
⁴² the oligosaccharyltransferases of family GT66, which transfer a pre-assembled oligosaccharide to Asp residues
⁴³ for protein N-glycosylation [4, 9]. Several sugar-diphospholipid-utilizing GTs are currently missing in the CAZy
⁴⁴ database, and here we classify new sugar-diphospholipid-utilizing GTs from four major functional classes that
⁴⁵ are all involved in the synthesis of bacterial cell wall polysaccharides.

⁴⁶ The first of these four functional classes corresponds to the peptidoglycan polymerases, SEDS (shape, elon-
⁴⁷ gation, division and sporulation) proteins. These proteins polymerize peptidoglycan in complex with class B
⁴⁸ penicillin-binding proteins [10]. Several 3-D structures of SEDS proteins have been determined, and they harbor
⁴⁹ 10 transmembrane helices and one long extracellular loop [11, 12, 13]. This loop contains an Asp residue, which
⁵⁰ has been shown to be essential for SEDS function [11, 14].

The enzymes in the next two functional classes, bacterial polysaccharide polymerases (BP-Pol, also known as Wzy) and O-antigen ligase (O-Lig, also known as WaaL) are involved in the synthesis of lipopolysaccharides (LPS). LPS are polysaccharides on the membrane of Gram-negative bacteria, and consist of the highly diverse O-antigen attached to the Lipid A-core oligosaccharide located in the outer membrane [15]. The structure of the O-antigen determines the O-serotype of the bacteria. Most LPS structures are produced via the so-called Wzx/Wzy-dependent pathway [16, 17], for which the genes are located in a specific gene cluster [16]. In this pathway, BP-Pol catalyzes the polymerization of pre-assembled oligosaccharides attached to Und-PP. Little is known about the activity of BP-Pols. Firstly, because they are difficult to express heterologously, and to date, only one study has demonstrated the activity of O-Pol *in vitro* [18] and no experimentally determined 3-D structure is available. Secondly, because the sequences of BP-Pols are highly diverse with a sequence identity as low as 16% for different serotypes of the same species [16], it is difficult to identify conserved residues. However, several studies have identified BP-Pols in the gene clusters of various species, paving the way for analyzing BP-Pol sequences across a large range of taxonomic origin (see below). These include some Gram-negative bacteria which also employ the Wzx/Wzy-dependent pathway to produce capsular polysaccharides, including *Streptococcus pneumoniae* [19]. The third functional class, O-Lig catalyzes the final step in the synthesis of LPS; the ligation of the newly synthesized polymer (O-antigen) onto Lipid A-core oligosaccharide [20]. A structure of O-Lig in complex with Und-PP has been reported, which showed a fold with 12 transmembrane helices and a long periplasmic loop containing several conserved residues; two Arg which bind to the phosphates of Und-PP and a His which is proposed to activate the acceptor [21].

The enzymes present in the fourth functional class, the enterobacterial common antigen polymerases (ECA-Pol, also known as WzyE) are involved in the synthesis of enterobacterial common antigen (ECA). In addition to the O-antigen, ECA is a specific polysaccharide that occurs on the cell surface in members of the Enterobacterales order. ECA consists of repeating units of N-acetylglucosamine, N-acetyl-D-mannosaminuronic acid and 4-acetamido-4,6-dideoxy-D-galactose [22]. ECA is also produced via the Wzy/Wzx-dependent pathway, where ECA-Pol performs the equivalent reaction to the BP-Pols.

Structurally, the sugar-diphospholipid-utilizing GTs have an overall GT-C fold common to other integral membrane GTs, which is different from the globular nucleotide-sugar-utilizing GTs; GT-A and GT-B [4]. GT-C enzymes have a number of transmembrane helices that varies from 8 to 14 [4, 23]. Alexander and Locher recently suggested two subgroups of GT-C glycosyltransferases, GT-C_A and GT-C_B [23], where O-Lig and SEDS make up GT-C_B [23]. As no structures have been published of ECA-Pol and BP-Pols, these have not been assigned to a structural subgroup.

We have identified 17 new GT families covering a large number of the sugar-diphospholipid-utilizing GTs, by detailed analysis of the primary sequence of SEDS proteins, ECA-Pols, BP-Pols and O-Ligs. In addition, we examined how sequence diversity correlates with the diversity of the transferred oligosaccharides and with the stereochemical outcome of the glycosyl transfer reaction. The analysis also revealed that the new GT families organize in three clans across the functional classes suggestive of common ancestry. Despite of poor sequence alignments we manage to identify conserved potentially critical amino acids common within the clans.

2 Results

2.1 Peptidoglycan Polymerases

For building the CAZy family of SEDS proteins, we used four characterized proteins as seed sequences: the proteins with PDB IDs 6BAR [11], 8TJ3 [13] and 8BH1 [12], and the protein with GenBank accession CAB15838.1 [24]. Family GTxx1 was created and initially populated by using BLAST against GenBank, and subsequently by searching against GenBank with an HMM built from the retrieved sequences. GTxx1 is a very large family currently counting over 57,200 GenBank members in the CAZy database with a pairwise sequence identity of 19% over 221 residues for the most distant members.

The taxonomic distribution of family GTxx1 follows what was reported in [14], namely that this protein family is present in all bacteria except for Mycoplasma. It is present in most but not all planctomycetes.

For SEDS proteins, the glycosyl donor for the polymerization reaction is Lipid II (Und-PP-muropeptide, an activated disaccharide carrying a pentapeptide), where the undecaprenyl diphosphate is α -linked. The carbohydrate repeat unit of peptidoglycan being β -linked, the glycosyl transfer reaction thus inverts the stereochemistry of the anomeric carbon involved in the newly formed glycosidic bond.

2.2 Enterobacterial common antigen polymerases

The ECA-Pol which was studied in [25] was used as seed sequence for building the ECA-Pol family. Although the CAZy database only lists GenBank entries [26], we decided to build our multiple sequence alignments (MSAs) with sequences from the NCBI non-redundant database in order to capture more diversity. An ECA-Pol sequence library was thus constructed from the seed sequence using BLAST against the non-redundant database of the

107 NCBI. The ECA-Pols were assigned to a single new CAZy family, GTxx2. To date this new family contains
108 over 4800 GenBank members with sequence identity greater than 38% over 414 residues, consistent with the
109 conservation of acceptor, donor and product of the reaction.

110 As expected from their taxonomy-based designation, the ECA-Pol family (GTxx2) essentially contains se-
111 quences from the Enterobacterales order but also a few members of the Pasteurellales, suggesting that ECA-Pols
112 of the latter were acquired by horizontal gene transfer.

113 The ECA-Pol family uses a retaining mechanism, since the substrate repeat unit is axially linked to Und-PP
114 and also axially linked in the final polymer.

115 2.3 O-antigen ligases

116 With the aim of including the O-Ligs in the CAZy database, we collected 37 O-Lig sequences (Supplemen-
117 tary Table 1) and constructed a sequence library from these seed sequences using BLAST against the NCBI
118 non-redundant database. A phylogenetic tree of the sequence library revealed four distantly related clades
119 (Supplementary Fig. 1). The O-Ligs were included into one new CAZy family, GTxx3 with more than 16,700
120 members distributed in the four subfamilies.

121 The greater diversity of the GTxx3 O-Ligs compared to the GTxx1 peptidoglycan polymerases and GTxx2
122 ECA-Pol appears in the form of the four divergent clades in the O-Lig phylogenetic tree (Supplementary Fig.
123 1). We hypothesize that this increased diversity originates from the extensive donor and moderate acceptor
124 variability of O-Ligs [15]. Taxonomically, the GTxx3 O-Lig family is present in most bacteria, including both
125 Gram-negatives and Gram-positive bacteria. The reaction performed by O-Ligs involves an inversion of the
126 stereochemistry of the anomeric carbon since the sugar donor is axially bound to Und-PP and the reaction
127 product is equatorially bound to Lipid A [20].

128 A recently discovered O-Lig, WadA, is bimodular with a GTxx3 domain appended to a globular glyco-
129 syltransferase domain of family GT25, which adds the last sugar to the oligosaccharide core [27]. We have
130 constructed a tree with representative WadA homologs from the GTxx3 family (Supplementary Fig. 2) and
131 observe that most of the sequences appended to a GT25 domain form one clade in the tree, except for a few
132 outliers. This suggests a coupled action of the GT25 and of the GTxx3 at least for the bimodular O-ligs and
133 possibly for the entire family. The bimodular WadA O-Lig is observed in five genera including Mesorhizobium
134 and Brucella.

135 2.4 Other bacterial polysaccharide polymerases

136 The fourth functional subgroup of GT-C_B are the BP-Pols. As previously mentioned, there is only one ex-
137 perimentally characterized BP-Pol [18], but several studies have identified BP-Pols from the polysaccharide
138 gene clusters, and we decided to build our families based on these published reports. We thus collected 363
139 predicted BP-Pol sequences from seven studies for various species, both Gram-negatives and Gram-positives: *Es-*
140 *cherichia coli* [28], *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri* [29], *Salmonella enterica* [30], *Yersinia*
141 *pseudotuberculosis*, *Yersinia similis* [31], *Pseudomonas aeruginosa* [16], *Acinetobacter baumanii*, *Acinetobacter*
142 *nosocomialis* [32] and *Streptococcus pneumoniae* [19] (Supplementary Table 2).

143 In contrast to ECA-Pols, the donors as well as the acceptors of BP-Pols are highly variable. Others have
144 reported an exceptional sequence diversity of BP-Pols even within the same species [16]. We also found that
145 the sequences of BP-Pols are extremely diverse, and global alignments failed to reveal any conserved residue
146 due to both sequence diversity and to the difficulty in aligning proteins with multiple and variable numbers of
147 transmembrane helices. It was therefore not possible to build a single family that could capture the diversity
148 of BP-Pols.

149 In order to group BP-Pols into similarity clusters that we could include as families in the CAZy database, we
150 first built a sequence library by running BLAST against the NCBI non-redundant database for each of the 365
151 BP-Pol seeds. Clustering of the BP-Pols proved challenging. A phylogenetic analysis was not possible because
152 of their great diversity, and a sequence similarity network (SSN) analysis alone would either result in very small
153 clusters (using a strict threshold) or larger clusters that were linked because of insignificant relatedness (using
154 a loose threshold).

155 Instead, we used a combination of SSN and HMM comparisons: First, we used an SSN with a strict threshold
156 which would allow us to build good MSAs for the resulting clusters. This resulted in 204 clusters (Fig. 1a).
157 Next, we created an HMM profile of each SSN cluster and compared the HMMs by all-vs-all pairwise HHblits,
158 a program that aligns two HMMs and calculates a similarity score [33]. We then combined the SSN clusters
159 into “superclusters” in a network analysis based on the HHblits scores (Fig. 1b), resulting in 28 superclusters of
160 varying sizes and 86 singleton clusters. Interestingly, the BP-Pols clustered across taxonomy, and even BP-Pols
161 from Gram-positive and Gram-negative bacteria clustered together. The 14 largest superclusters define new
162 GT families in the CAZy database (GTxx4-GTx17) with a number of members ranging from 159 to 5,979 at

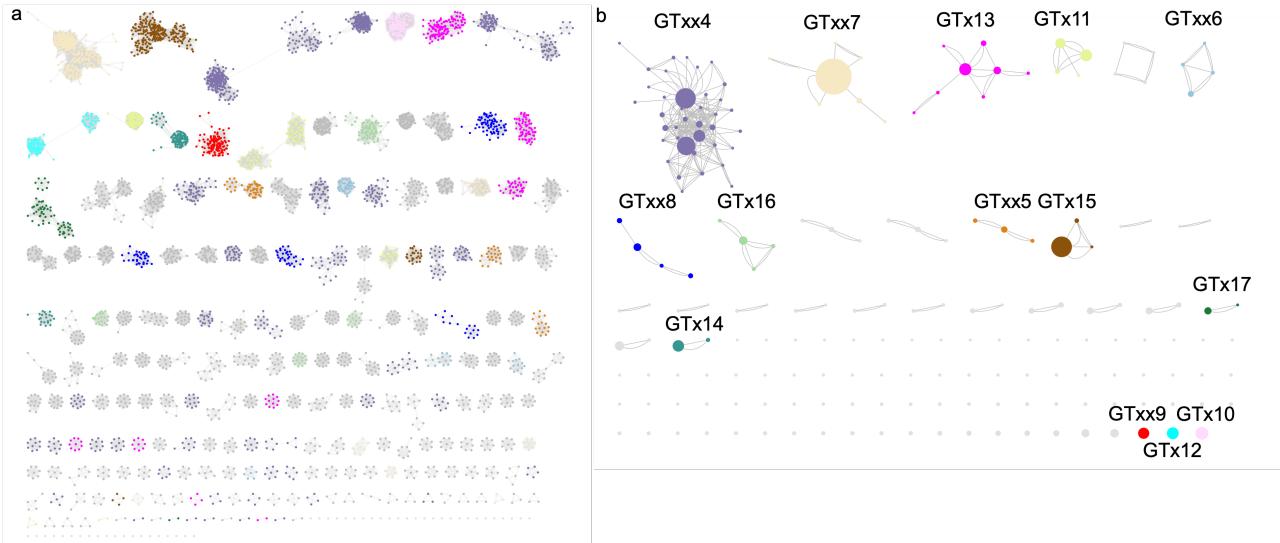


Figure 1: Clustering of BP-Pol sequences. a) SSN network with nodes representing proteins and edges representing pairwise alignment bit scores. b) HHblits network with nodes presenting SSN clusters and edges representing HHblits scores. The resulting clusters are referred to as “superclusters”. There are two edges between nodes, when the HHblits score is above the threshold in both directions. The size of the nodes represents the number of members in the SSN cluster. The 14 largest superclusters (>150 GenBank members) defined CAZy families GTxx4 - GTx17. In both a and b, the SSN clusters are coloured according to which supercluster they belong to.

the time of submission. Only 150 of the 363 original seeds are included in the new families. We thus expect that many more BP-Pol families will be created in the future, as the amount and diversity of data increase.

All of the BP-Pol families are present in a wide taxonomic range, and outside of the taxonomic orders of the original seeds. Several of the families contain members from both Gram-positive and Gram-negative bacteria, for example GTxx4, GTx12, and GTx16.

As a way of evaluating our families, we performed structural superimpositions of AlphaFold models of distantly related members of each family. As an example, five distantly related members of GTxx4 are shown in Supplementary Fig. 4. The sequence identity between these members is relatively low (between 21.4 and 24.3%). Yet, they still produce a meaningful superimposition, and notably, the conserved residues are oriented very similarly.

2.5 Analyzing the sugars transferred by bacterial polysaccharide polymerases

Next, we investigated how the BP-Pol families relate to the structures of the transferred oligosaccharide repeat units. We retrieved the serotype-specific sugar structures, which were reported in the review papers [34, 29, 30, 31, 16, 32, 19]. Additionally, nine sugar structures were included, which were published after the review papers [35, 36, 37, 38, 39]. Out of the 150 BP-Pol seed sequences that were included in the new CAZy families, we matched 131 with a sugar structure. The repeat units are oligosaccharides with 3-7 monomers within the backbone, often with branches. In most of the cases, the bond which is formed by the polymerase has been identified in the review papers based on the other GTs in the gene cluster which assemble the repeat units.

Having retrieved the sugar structures, we first analyzed the stereochemistry of the bond catalyzed by the polymerase. As mentioned above, the stereochemical mechanism (inverting or retaining) is usually well conserved in the CAZy GT families. The repeat unit structures are always axially linked (α for D-sugars and β for L-sugars) to the Und-PP moiety before polymerization. There are two possible mechanisms for the BP-Pol-catalyzed polymerization reaction, either retaining or inverting the axial configuration. Thus, if the bond formed by the polymerase is axial, the mechanism is retaining and if the bond formed by the polymerase is equatorial, the mechanism is inverting.

We found that the stereochemical outcome of BP-Pols appears well conserved within the new BP-Pol CAZy families and varies from one family to another (Fig. 2). There is only one exception; in family GTxx8, the polymerase linkages are all equatorial except for the O-antigen in *Pseudomonas aeruginosa* O4, where it is axial. It is possible that there could be an error in the chemical structure or that the serotype designation was incorrect or that the *P. aeruginosa* O4 polymerase constitutes an exception.

Next, we investigated whether there was a correlation between the structures of the transferred sugars and the sequence similarity of the BP-Pols. We created phylogenetic trees of the BP-Pols in each family and

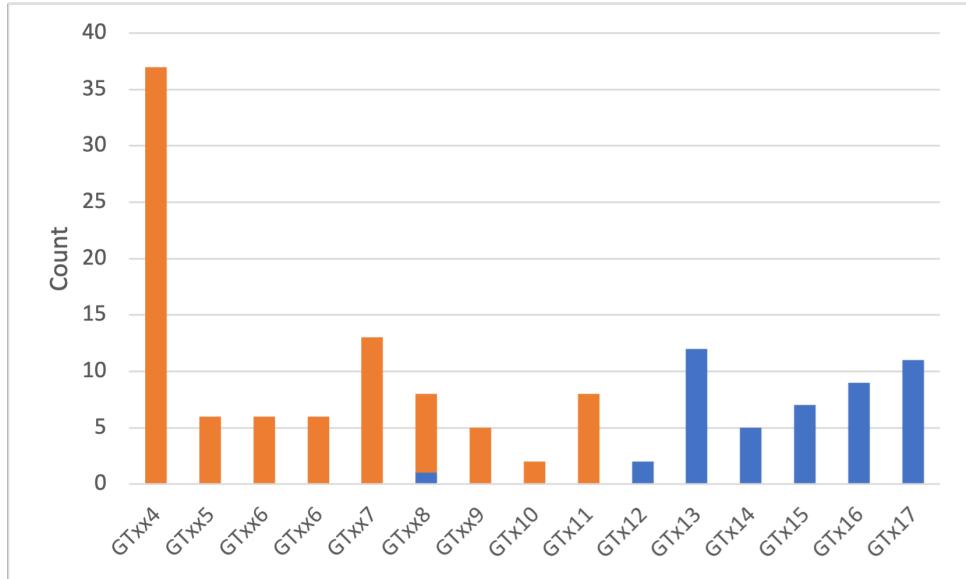


Figure 2: Conservation of stereochemical outcome of the reaction catalyzed in the various BP-Pol families. Equatorial bonds are shown in orange, implying an inverting mechanism. Axial bonds are shown in blue, implying a retaining mechanism.

195 visualized them with the corresponding transferred repeat units. We observe that the sugars within each family
 196 show similarity and this similarity appears to correlate with the structure of the tree (Fig. 3, Supplementary
 197 Fig. 4). The ends of the repeat units, ie. the subsite moieties immediately upstream (+1) and downstream (-1)
 198 of the newly created bond (Fig. 4) seem to be most conserved whereas more variability occurs in the middle
 199 part. We hypothesize that the +1 and -1 subsites are the moieties most important for recognition by the active
 200 site of the BP-Pol.

201 We observe examples of BP-Pols from distant taxonomic origin that cluster in the same CAZy family and
 202 have highly similar sugars. For example, *Escherichia coli* O178 and *Streptococcus pneumoniae* 47A in GTxx7
 203 transfer sugars with almost identical backbones, suggestive of horizontal gene transfer. There is only a slight
 204 variance in the middle of the repeat unit. This suggests that there is less constraints on the central part of the
 205 repeat unit than on the extremities that define the donor and the acceptor.

206 We next attempted to quantify the correlation between BP-Pol sequence and carbohydrate structure. For
 207 this we developed an original pairwise oligosaccharide similarity score. In our scoring scheme, the similarity of
 208 two glycans is estimated by examining the -1 and +1 subsites, as we expect that these are the moieties most
 209 fitting the active site of the BP-Pol (Fig. 4). The minimum match between two oligosaccharides corresponds to
 210 identical moieties at both subsites -1 and +1, which yields a score of 2. Thereafter, the score increases by one
 211 unit for each additional match at contiguous subsites, -2, -3, etc., and +2, +3, etc., up to a maximum value of
 212 7 subsites found for the glycans encountered in this study (for details see Methods).

213 Using our glycan similarity scoring system, we found a correlation between sugar similarity and polymerase
 214 sequence similarity (Fig. 5), supported by a preponderance of similarity scores appearing close to the score
 215 matrix diagonal and within each individual family.

216 2.6 Comparison of families

217 Others have previously reported sequence and structural similarity between SEDS, O-Lig and some BP-Pols
 218 [13, 23, 21, 14]. In order to investigate the relatedness of the new CAZy families, we compared the family
 219 HMMs by all-vs-all HHblits analyses [33] (Fig. 6). Strikingly, we observe that the retaining BP-Pol families
 220 cluster together on the heatmap along with the retaining ECA-Pols, while the inverting BP-Pols form two
 221 distinct groups, one of them containing the inverting O-Ligs. The background noise between some inverting
 222 and retaining enzymes is likely due to the general conservation of the successive transmembrane helices, which
 223 is altered in the GTxx4-GTxx5-GTxx6 subgroup due to their different architecture (see below). Peptidoglycan
 224 polymerases, GTxx1, segregate away from the other families.

225 In the CAZy database, clans have been defined for the glycoside hydrolases (GHs), which group together
 226 CAZy families with distant sequence similarity, similar fold, similar catalytic machinery and stereochemical
 227 outcome [40]. In extension of the report of the GT-CB class by Alexander and Locher [23], and based on the
 228 above-mentioned similarities between the new CAZy families, we can now define three clans within GT-C_B:
 229 GT-C_{B1} consisting of inverting BP-Pol families and O-Lig, GT-C_{B2} consisting of retaining BP-Pol families and

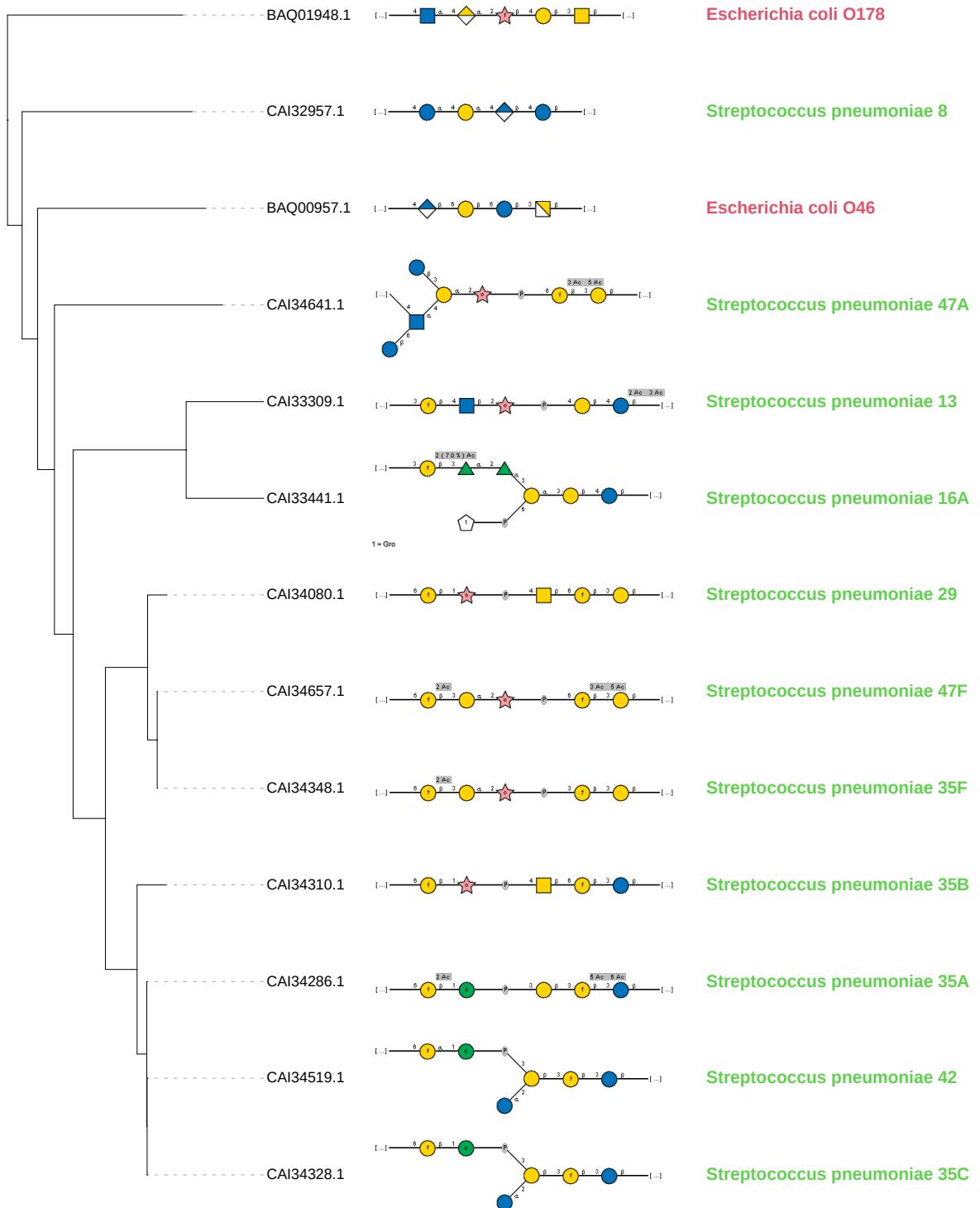


Figure 3: Similarity of transferred sugars by BP-Pols in GTxx7. The transferred repeat unit structures (in SNFG representation) are shown on a phylogenetic tree of BP-Pols in family GTxx7. There is an overall similarity between all the transferred sugars in the family and the similarity appears to correlate with the tree structure, ie. BP-Pol similarity. In particular, the ends of the repeat units (+1 and -1 subsites) appear to be often conserved, whereas there is more variety in the central region. Notably, the family contains BP-Pols from distant taxonomy which transfer similar sugars.

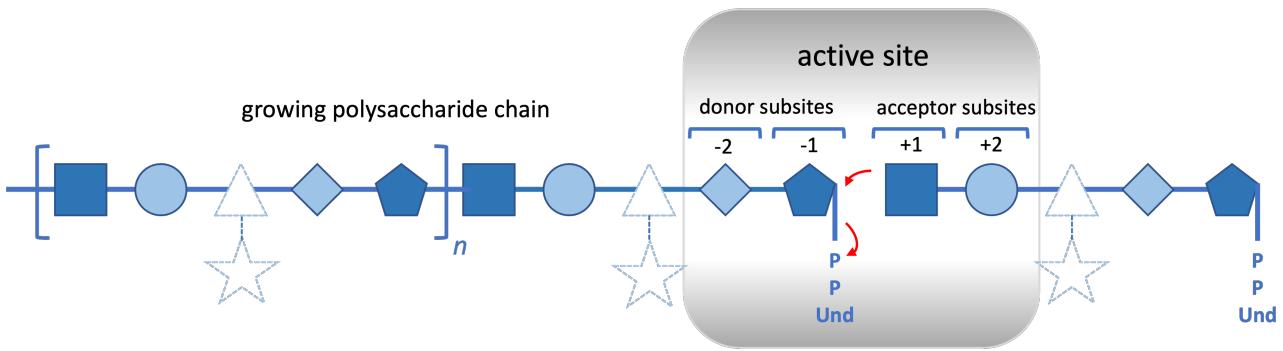


Figure 4: An idealized representation of a BP-Pol. The donor is the growing glycan chain activated by undecaprenyl pyrophosphate while the acceptor is a repeat unit monomer. The reaction is hypothesized to chiefly involve the sugar residues of the donor (subsites -2 and -1) and of the acceptor (subsites +1 and +2) that are proximal to the reaction center rather than residues and branches that are more distal (depicted in dashed lines). The reaction is represented by red arrows.

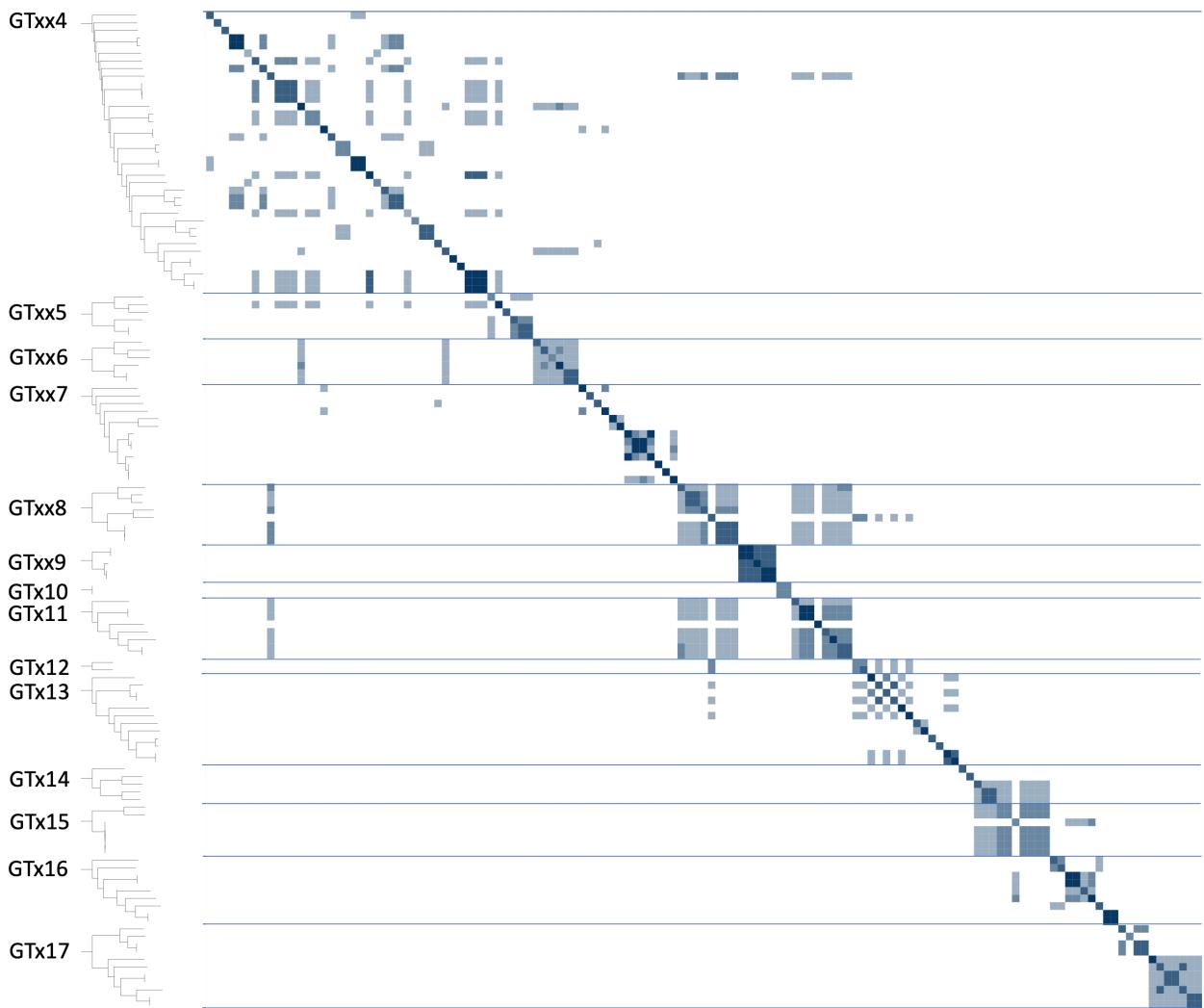


Figure 5: Glycan similarity of sugar repeat units polymerized by BP-Pols. All “seed” BP-Pols where the corresponding transferred oligosaccharide was known were included in the heatmap. A phylogenetic tree is shown for the polymerases in each CAZy family on the left. The glycan similarity scores are shown in a color scale of light blue (score value of 2 corresponding to identical matches at both -1 and +1 sites) to dark blue (score value of 5 corresponding to identical matches for at least three additional sequential positions). Blue lines separate the families.

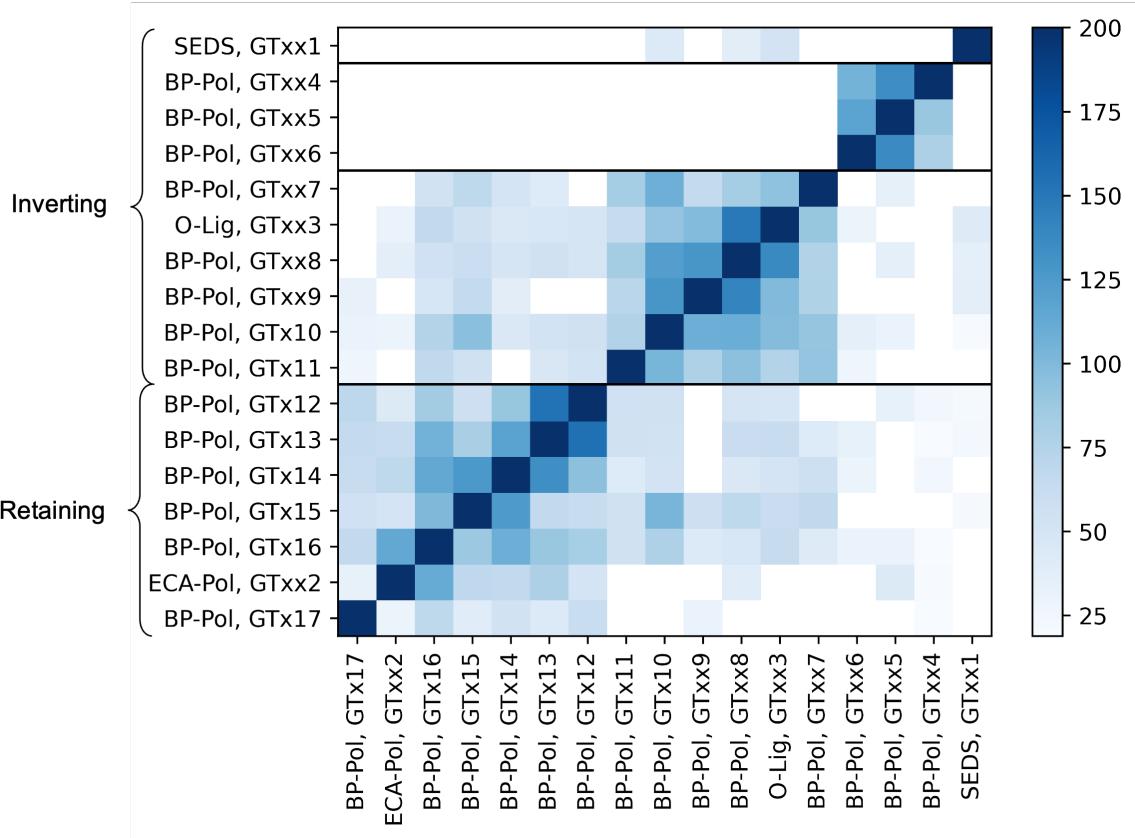


Figure 6: Heatmap of inter-family HHblits bit scores. The HHblits scores are shown on a color scale from white (low similarity score) to dark blue (high similarity score). The HHblits scores depend on the direction of the alignment, and therefore the heatmap is not symmetrical.

ECA-Pol, and GT-C_{B3} consisting of inverting BP-Pol families (Table 1). The families within each clan share residual, local, sequence similarity, insufficient to produce a multiple sequence alignment, but suggestive of common ancestry.

In the absence of a three-dimensional structure, and based solely on the number of transmembrane helices, we assigned clan GT-C_{B2} and GT-C_{B3} to the structural subclass GT-C_B of Alexander and Locher [23]. In addition, we also present in Table 1 the families of GT-C glycosyltransferases that have not yet been assigned to a structural class.

We then examined residue conservation and the general architecture of the enzymes in the clans. Based on the above mentioned pairwise HHblits analyses and structural superimpositions (Supplementary Figure 5-7), we tried to evaluate which architectural features and conserved residues are common within the clans. Indeed, there are some common features across most families. In all the families, all the conserved residues are located on the outer face of the membrane. Enzymes of clans GT-C_{B1} and GT-C_{B2} have a long extracellular loop close to the C-terminus (Fig. 7). In stark contrast, families GTxx4, GTxx5 and GTxx6 of clan GT-C_{B3} have an architecture completely different from that of the two other clans (Fig. 7), with the long loop located close to the N-terminus, and a conservation of one Asp, one His and two Arg residues.

As mentioned above, the structure of O-Lig in complex with Und-PP revealed several important residues; Arg-191 and Arg-265 which bind to the phosphate groups of Und-PP, and His-313 which is proposed to activate the acceptor [21]. The other families in GT-C_{B1} appear to have a similar pattern. All have 1-2 conserved Args, most of which are conserved in the HHblits alignment, and we hypothesize that they also play the role of binding to the diphosphate. Similarly, all families in the clan except for GTxx9 have either a conserved Asp or Glu, which align with the His-313 in O-Lig. We hypothesize that the Glu and Asp residues in the BP-Pols play the same role as the His-313 in O-Lig. As an example, the superimposition of the published O-Lig structure (7TPG) [21] and an AlphaFold model from one representative of the inverting BP-Pol family GTxx8 is shown in Fig. 9a. The superimposition produced an overall RMSD of 5.3 Å over 192 residues. Even with such a high RMSD, the two conserved Args are oriented very similarly, and the conserved His in O-Lig is placed in the same position as the conserved Glu in the BP-Pol.

In the retaining clan GT-C_{B2}, the pattern of conservation is different. Here, most of the families have 2-3 conserved Arg/Lys and 1-2 conserved Tyr. Interestingly, we observe that the ECA-Pol family GTxx2 shows high similarity with one of the BP-Pol families, GTx16. A superimposition of AlphaFold models from each

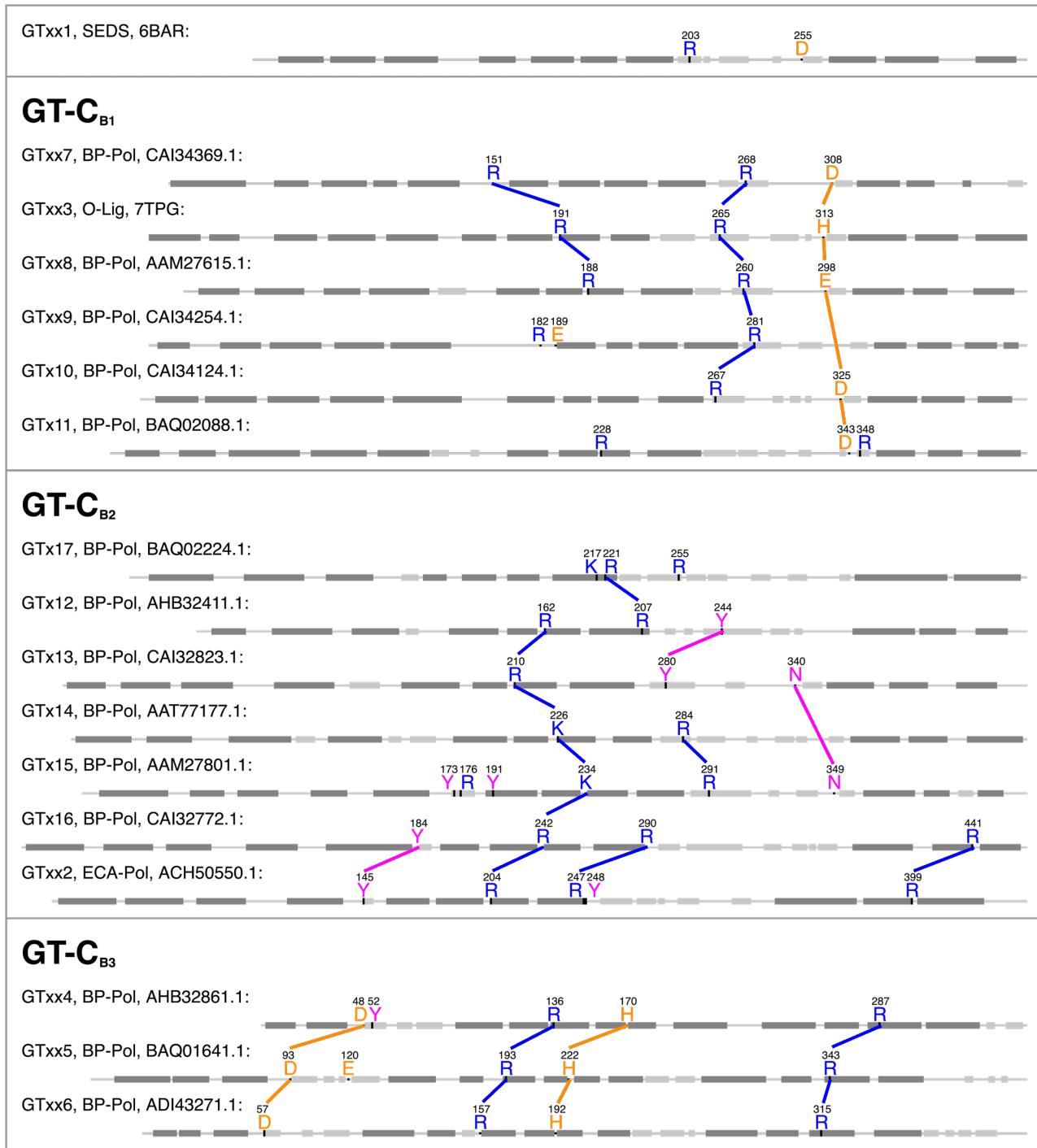


Figure 7: Equivalent conserved residues in the clans. The conserved residues of each of the new CAZy families are shown on sequences of representative family members. Lines are shown between conserved residues that align in HHblits alignments and that co-localize in structural superimpositions (Supplementary Fig. 5). Transmembrane helices are shown in dark gray boxes, non-transmembrane helices are shown in light gray boxes. The secondary structures were taken from the crystal structures for family GTxx1 and GTxx3 (6BAR and 7TPG respectively) and from AlphaFold models for all other families. The R210 in GTx13 is either K or R in the family. Conserved aliphatic residues are not shown.

Structural subclass Alexander & Locher	CAZy clan	CAZy families	Mechanism	Donor
GT-C _A (7 conserved TM helices)	-	GT53	Inverting	Lipid-P-monosaccharide
	-	GT83	Inverting	Lipid-P-monosaccharide
	-	GT39	Inverting	Lipid-P-monosaccharide
	-	GT57	Inverting	Lipid-P-monosaccharide
	-	GT66	Inverting	Lipid-PP-oligosaccharide
GT-C _B (10 conserved TM helices)	-	GTxx1	Inverting	Lipid-PP-oligosaccharide
	GT-C _{B1}	GTxx3, GTxx7, GTxx8, GTxx9, GTx10, GTx11	Inverting	Lipid-PP-oligosaccharide
	GT-C _{B2}	GTxx2, GTx12, GTx13, GTx14, GTx15, GTx16, GTx17	Retaining	Lipid-PP-oligosaccharide
	GT-C _{B3}	GTxx4, GTxx5, GTxx6	Inverting	Lipid-PP-oligosaccharide
-	-	GT22	Inverting	Lipid-P-monosaccharide
	-	GT50	Inverting	Lipid-P-monosaccharide
	-	GT58	Inverting	Lipid-P-monosaccharide
	-	GT59	Inverting	Lipid-P-monosaccharide

Table 1: Structural subclasses, clans and families of GT-C fold glycosyltransferases and relationships to mechanism and glycosyl donor.

family shows that the conserved residues are oriented very similarly, despite the low overall similarity (RMSD 5.4 Å over 360 residues) (Fig. 8b).

Although the peptidoglycan polymerase family, GTxx1 does not cluster in any of the three clans, it does display topographical similarity to clan GT-C_{B1}. In terms of architecture it also contains a long extracellular loop with a conserved Arg and the conserved and essential Asp residue [11]. The Asp residue is in a similar position as the Asp/Glu/His in the other families in clan GT-C_{B1}. We therefore hypothesize that this conserved Asp may play the role of activating the acceptor in clan GT-C_{B1} glycosyltransferases as the His in O-Lig [21].

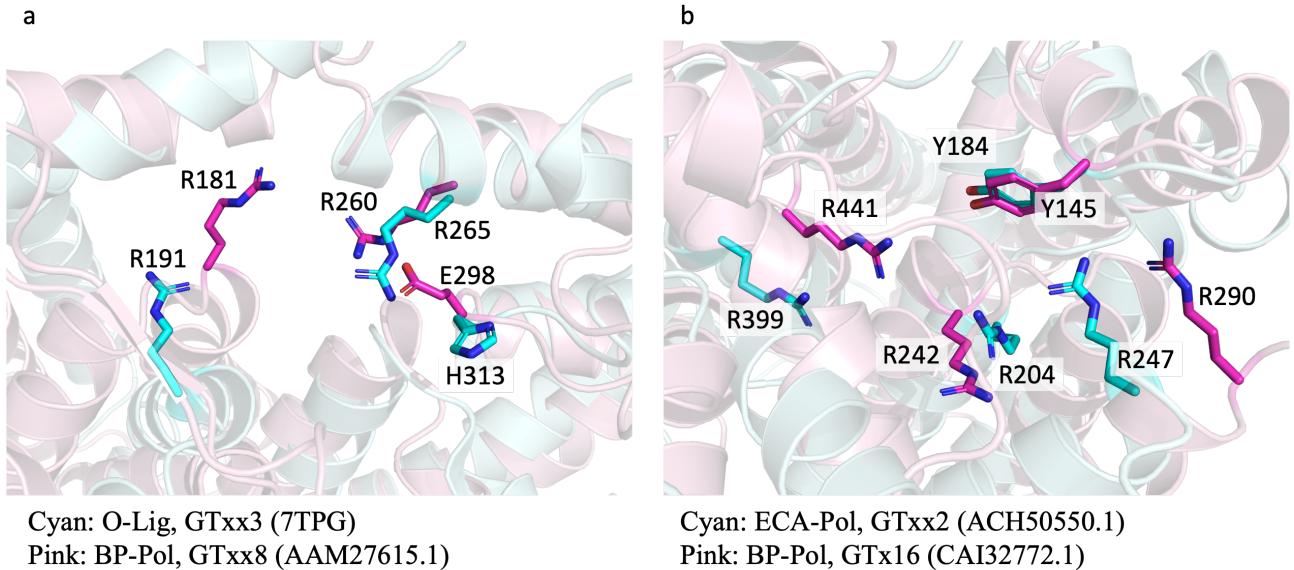


Figure 8: Structural superimposition of different families with conserved residues. a) O-Lig from GTxx3 (PDB 7TPG) and AlphaFold model of BP-Pol from GTxx8 (RMSD 5.3 Å over 192 residues). The conserved Glu in GTxx8 is aligning with the conserved His in GTxx3, which is proposed to activate the acceptor [21]. b) AlphaFold models of ECA-Pol from GTxx2 and BP-Pol from GTx16 (RMSD 5.4 Å over 360 residues). The conserved residues occupy similar positions.

266 3 Discussion

267 Here we have added 17 glycosyltransferase families (GTxx1 to GTx17) to the CAZy database bringing the
268 total of covered families from 116 to 133. In the CAZy database, families are built by aggregating similar
269 sequences around a biochemically characterized member. The known difficulties in the direct experimental
270 characterization of integral membrane GTs render this constraint impractical. To circumvent this problem, but
271 to remain connected to actual biochemistry, we decided to build our families around seed sequences for which
272 knowledge of the glycosidic bond formed could be deduced from examination of the polysaccharide product
273 from the literature.

274 To our knowledge, this is the first time that BP-Pols from different species have been successfully clustered.
275 Indeed, forming groups of BP-Pols has been very difficult previously because of their extreme diversity even
276 within strains of a single species [28], and, as a consequence, the knowledge on conserved and functional residues
277 has been very limited. By combining BP-Pols from a wide range of taxonomical origins and expanding with
278 the current sequence diversity, we were able to form larger families of similar polymerases from widely different
279 taxonomies, thereby revealing conserved residues that are most likely functionally important.

280 We observed that the O-Lig family (GTxx3) was present in many Gram-positive bacteria such as *Streptococcus*
281 *pneumoniae*. The covalent anchoring of CPS in Gram-negative bacteria is still poorly understood, although
282 it is found to be linked to peptidoglycan in some Gram-positive bacteria [17, 41]. Thus a hypothesis could be
283 that the GTxx3 members in *S. pneumoniae* are responsible for the ligation of CPS to the peptidoglycan layer
284 in these bacteria.

285 Because families are more robust when built with enough sequence diversity, many clusters of O-antigen
286 polymerases were judged too small to build meaningful CAZy families. Additional polymerase families are thus
287 expected in the future with the accumulation of sequence data. For instance the small cluster that contains
288 47% identical BP-Pols from *E. coli* (GenBank BAQ01516.1) and *A. baumanii* (GenBank AHB32586.1) only
289 contains eight sequences and will remain unclassified until enough sequence diversity has accumulated. This
290 arbitrary decision comes from the need to devise a classification that can withstand a massive increase in the
291 number of sequences without the need to constantly revise the content of the families.

292 Moreover, we observe that the sequence diversity within the families we have built is minimal for peptido-
293 glycan polymerases (GTxx1), and then increases gradually from ECA-Pols (GTxx2) to O-Ligs (GTxx3) and is
294 maximal for BP-Pols (GTxx4-GTx17). We hypothesize that sequence diversity reflects the donor and acceptor
295 diversity in each family since the latter increases accordingly.

296 It has been observed that for classical GT-A and GT-B fold glycosyltransferases, the catalytic mechanism is
297 conserved within a family, but families with the same fold can have different mechanisms, possibly because the
298 stereochemical outcome of the glycosyl transfer reaction is essentially dictated by the precise positioning and
299 activation of the acceptor above (S_N2) or below (S_{Ni}) the sugar ring of the donor [4]. Very occasionally, retaining
300 glycosyltransferases have been shown to operate via a double displacement mechanism that involves Asp/Glu
301 residues to form a glycosyl enzyme intermediate and to activate the acceptor that attacks this intermediate [42].
302 The families defined here display globally similar GT-C folds, and they also show conservation of the catalytic
303 mechanism with about half of the families retaining and the other half inverting the anomeric configuration of
304 the donor, suggesting that the outcome of the reaction catalyzed by GT-C glycosyltransferases is also dictated
305 by the positioning of the acceptor with respect to the sugar plane of the acceptor. In turn this also suggests
306 that retaining BP-Pols also operate by an S_{Ni} mechanism rather than by the formation of a glycosyl enzyme
307 intermediate. This hypothesis is supported by the lack of invariant Asp or Glu residues which could be involved in
308 the formation and subsequent breakdown of a glycosyl enzyme intermediate in the retaining families GTxx2 and
309 GTx12-GTx17. Additionally, the S_{Ni} mechanism may provide protection against the interception of a glycosyl
310 enzyme intermediate by a water molecule resulting in an undesirable hydrolysis reaction and termination of the
311 polysaccharide elongation.

312 The wealth of structural data of GT-C glycosyltransferases now permits a deeper evaluation of the intrinsic
313 properties of this large class of enzymes. Alexander and Locher have recently evaluated the structural similarities
314 between GT-C fold glycosyltransferases and have divided them in two fold subclasses [23]. The GT families
315 that we describe here significantly expand the GT-C class in the CAZy database (www.cazy.org) and allow to
316 combine the structural classes with mechanistic information. Lairson *et al.* have proposed the subdivision of
317 GT-A and GT-B fold glycosyltransferases in clans that integrate the stereochemical outcome of the reaction [4].
318 Here we also note the conservation of the stereochemistry in the families of BP-Pols and we thus propose to
319 group them into three clans which share the same fold, residual sequence conservation and the same catalytic
320 mechanism (Table 1). As more families of BP-Pols emerge, these three clans will likely grow. Table 1 shows
321 the three clans we defined here and how they relate to the structural classes defined by Alexander and Locher.
322 Of note are families GTxx4, GTxx5, and GTxx6 which do not bear any similarity, even distant, with the GT
323 families of the other two clans. These three families also stand out by the location in the sequence of the long
324 loop that harbors the catalytic site in the other GT-C families. In absence of relics of sequence relatedness to
325 the other families, GTxx4, GTxx5 and GTxx6 were assigned to clan GT-C_{B3}. With 10 transmembrane helices,

326 it is tempting to suggest that this clan may belong to the fold subclass GT-C_B of Alexander and Locher.

327 The analysis presented here shows that not only the stereochemistry of the glycosyl transfer is conserved
328 in the BP-Pol families, but our development of an original method to estimate glycan similarity also reveals a
329 certain degree of structural similarity of the oligosaccharide repeat units, suggesting that the latter constitutes
330 a significant evolutionary constraint applying to the sequence and structure of BP-Pols. A closer inspection
331 of the oligosaccharide repeat units within the families further reveals that the carbohydrates that appear the
332 most constrained are the carbohydrates located (i) at the non-reducing end of the acceptor and (ii) close to the
333 undecaprenyl pyrophosphate of the donor, i.e. the residues closest to the reaction center (Fig. 4). By contrast,
334 residues away from the two extremities engaged in the polymerization reaction appear more variable, and can
335 tolerate insertions/deletions or the presence of flexible residues such as linear glycerol or ribitol, with or without
336 or the presence of a phosphodiester bond.

337 The version of the glycan similarity score presented here was inspired in part by observed structural simi-
338 larities in different O-antigen repeat units assembled by very similar BP-Pols [16]. The repeat-unit comparison
339 involves a translation of glycan IUPAC nomenclature to a reduced alphabet of terms representing only backbone
340 configuration, i.e., ignoring chemical modifications and sidechains. Furthermore, a positive similarity score re-
341 quires an entire identical match of all backbone elements at both donor and acceptor positions (-1 and +1 sites
342 in Fig. 4, respectively). Despite these simplifications, the similarity score reveals, with exceptions, an overall
343 greater intra- rather than inter-family oligosaccharide similarity (Fig 5). These limitations will be addressed at
344 a later stage (G.P. Gippert, in preparation).

345 We have next looked at the distribution of the new GT families in genomes, and particularly the families
346 of BP-Pols. This uncovers broadly different schemes, with some bacteria having only one polymerase (and
347 therefore only able to produce a single polysaccharide) while others having several, and sometimes more than
348 5, an observation in agreement with the report that *Bacteroides fragilis* produces no less than 8 different
349 polysaccharides from distinct genomic loci [43]. The multiplicity of polysaccharide biosynthesis loci in some
350 genomes makes it sometimes difficult to assign a particular polysaccharide structure to a particular biosynthesis
351 operon.

352 As already shown in other occasions, the sequence-based classification of carbohydrate-active enzymes of
353 the CAZy database has predictive power. The case of the GT families described here supports this view as
354 the invariant residues in the families not only co-localize in the same area of the three-dimensional structures
355 (whether actual or AlphaFold-predicted), but also correspond to the residues found essential for function in
356 the families where this has been studied experimentally. The families described herein also show mechanistic
357 conservation and thus the stereochemistry of glycosyl transfer can be predicted. Finally, the observed similarity
358 in oligosaccharide repeat units that accompanies sequence similarity has also predictive power and paves the
359 way to the future possibility of *in silico* serotyping based on DNA sequence.

360 4 Methods

361 4.1 Alignment-based Clustering (Aclust)

362 Phylogenetic trees were generated using an in-house tool called Aclust (G.P.Gippert, manuscript in preparation).
363 Source code may be obtained via GitHub at <https://github.com/GarryGippert/Aclust>. Aclust employs a
364 hierarchical clustering algorithm comprising the following steps. (1) A distance matrix is computed from all-vs-
365 all pairwise local pairwise sequence alignments [44], or from a multiple sequence alignment provided by MAFFT
366 [45]. The distance calculation is based on a variation of Scoredist [46] where distance values are normalized
367 to the shorter pairwise sequence length rather than to pairwise alignment length. (2) The distance matrix is
368 embedded into orthogonal coordinates using metric matrix distance geometry [47], and (3) a bifurcating tree is
369 computed using nearest-neighbor joining and centroid averaging in the orthogonal coordinate space. The last
370 centroid created in this process is defined as the root node. (4) Beginning with the root node of the initial
371 tree, each left and right subtree constitutes disjoint subsets of the original sequence pool, which are reembedded
372 and rejoined separately (i.e., steps 2 and 3 repeated for each subset), and the process repeated recursively —
373 having the effect of gradually reducing deleterious effects on tree topology arising from “long” distances between
374 unrelated proteins.

375 4.2 Building the peptidoglycan polymerase family (GTxx1)

376 The peptidoglycan polymerase family, GTxx1, was built by using Blastp from BLAST+ 2.12.0+ [48] with the se-
377 quences of the characterized SEDS proteins (PDB 6BAR, 8TJ3 and 8BH1 and GenBank accession CAB15838.1)
378 against GenBank with a threshold of approximately 30% to retrieve the family members. Next, an MSA was
379 generated with MAFFT v7.508 using the L-INS-i strategy [45], and an HMM model was built with hmmbuild of
380 HMMER 3.3.2 [49]. The family was further populated using hmmsearch from HMMER 3.2.2 against GenBank.

4.3 Building the Enterobacterial common antigen polymerases family (GTxx2)

A sequence library of ECA-Pols was constructed by using Blastp with the seed sequence (GenBank accession AAC76800.1) against the NCBI non-redundant database version 61 with an E-value threshold of 1e-60. The hits were redundancy reduced using CD-HIT 4.8.1 [50] with a threshold of 99%. The redundancy-reduced pool of ECA-Pol sequences was clustered using our in-house tool Aclust (see above), and the tree showed one large clade and a few outliers. All the sequences in the large clade were used to build an MSA using MAFFT v7.508 with the L-INS-i strategy [45]. An HMM was built based on this MSA using hmmbuild of HMMER 3.3.2 [49]. The family GTxx3 was built in CAZy and populated using Blastp against GenBank with an approximate threshold of 30% and hmmsearch against GenBank.

4.4 Building the O-antigen ligase family (GTxx3)

37 O-Lig sequences were selected from literature (Supplementary Table 1) and expanded using Blastp against the NCBI non-redundant database with an E-value cut-off of 1e-60. Redundancy reduction was performed on the resulting sequence pool using CD-HIT with a threshold of 99%, resulting in a pool of 1,402 sequences. A phylogenetic tree of the pool of O-Lig sequences was generated using Aclust (see above), which showed deep clefts between main branches, and branches with sufficient internal diversity (Supplementary Figure 2). Based on these results, four subfamilies were determined. An MSA was built for the family as well as for the subfamilies with MAFFT v7.508 using the L-INS-i strategy. HMMs were built based on the MSAs using the hmmbuild of HMMER 3.3.2 [49]. The family was populated using Blastp against GenBank using an approximate threshold of 30% identity with the seed sequences and using hmmsearch with the family and subfamily HMMs.

4.5 Building the Bacterial polysaccharide polymerase families (GTxx4-GTx17)

363 BP-Pol sequences were retrieved from review papers on biosynthesis of O-antigens and capsular polysaccharides in different species: *Escherichia coli* [28], *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri* [29], *Salmonella enterica* [30], *Yersinia pseudotuberculosis*, *Yersinia similis* [31], *Pseudomonas aeruginosa* [16], *Acinetobacter baumanii*, *Acinetobacter nosocomialis* [32] and *Streptococcus pneumoniae* [19] (complete list in Supplementary Table 2). The BP-Pols for *A. baumannii* O7 and O16 were omitted, because of uncertainty of their serotypes [32]. The BP-Pol from *P. aeruginosa* O15 was also omitted, because it has been shown that this BP-Pol is inactivated and that the O-antigen is synthesized via the ABC-dependent pathway rather than the Wzx/Wzy-dependent pathway [51].

The sequence library was expanded using Blastp for each seed sequence against the NCBI non-redundant database with an E-value threshold of 1e-15. Redundancy reduction was performed using CD-HIT with a threshold of 95% identity.

To find clusters of BP-Pol sequences that were large enough to create a CAZy family, we developed a clustering method consisting of two steps. First, in order to make a sequence similarity network (SSN), all-vs-all pairwise local alignments of the BP-Pol sequence pool were performed using Blastp from BLAST+ 2.12.0+. A series of networks were built using different bit score thresholds. The members of the resulting SSN clusters were identified using NetworkX [52] and MSAs of the members were built with MAFFT v7.508 using the L-INS-i strategy. The MSAs were inspected using Jalview [53], and a bit score threshold of 110 was selected, as it was the lowest score for which the SSN clusters had adequate sequence conservation (approximately 15 conserved residues).

HMMs were then built for each SSN cluster using hmmbuild of HMMER 3.3.2, and the HMMs were compared using HHblits 3.3.0 [54]. A series of HHblits networks were built using different HHblits score thresholds. Again, the members of the resulting “superclusters” were identified using NetworkX and MSAs of the members were built with MAFFT v7.508 using the L-INS-i strategy. A bit score threshold of 160 was selected as it resulted in “superclusters” with adequate diversity for building CAZy families (approximately 5 conserved residues). CAZy families were created for the 14 largest superclusters and populated with sequences present in GenBank by a combination of Blastp with the seed sequences and hmmsearch. The networks were visualized with Cytoscape [55].

4.6 Analysis of sugar repeat unit structures

In order to analyze the relation between BP-Pol sequence and structure of the transferred repeat unit, we retrieved the repeat unit structures for the serotypes for the BP-Pols that were included in the new CAZy families. The repeat unit structures were retrieved from the same review papers from which we got the BP-Pol sequences [32, 19, 31, 30, 29, 16], except for the sugars for *E. coli*, where the sugar structures have been reported elsewhere [34]. Nine additional repeat unit structures were included for *S. pneumoniae*, which were published after the review paper; serotypes 16A [35], 33A [36], 33C and 33D [37], 35C and 35F [38], 42 and 47F [56] and 47A [57]. For *Y. pseudotuberculosis* O3 and *S. pneumoniae* 33B, we used the revised structure from [39]

and [37] respectively. *Pseudomonas aeruginosa* O2 and O16 contain two BP-Pol genes; one BP-Pol localized in the O-antigen biosynthesis cluster, which polymerizes the sugar repeat units with an α bond and one BP-Pol localized outside the biosynthesis cluster which polymerizes the repeat units with a β bond [58]. Since the BP-Pol reported in [16] are the BP-Pols from the O-antigen cluster, we report the sugar structure with the α bond.

The linkages formed by the polymerase has been determined in all of these papers, except for a few cases. This determination is based on the initial GT transferring specific monosaccharides, and sometimes also based on other GTs in the gene cluster. The cases where the polymerase linkage has not been determined unambiguously are *E. coli* O166, O78, O152, O81, O83, O11, O112ab, O167, O187, O142, O117, O107, O185, O42, O28ac, O28ab, for which there were two or more possible linkages. For *S. pneumoniae* 33A, we determined the polymerase linkage based on the gene cluster having the initial transferase WchA, which transfers a glucose [32]. *S. pneumoniae* 47A has WcjG as the initial transferase, which transfers Galp or Galf [32]. Since the repeat unit contains both Gal and Galp, we could not determine the polymerase linkage unambiguously. However, the repeat unit is very similar to other repeat units in the family (most similar to that of *S. pneumoniae* 13), and we proposed the equivalent polymerase linkage.

The CSDB database (<http://csdb.glycoscience.ru>) [59] was used to retrieve literature, SNFG image representations and linear sugar strings of the repeat unit structures. Phylogenetic trees for BP-Pol families with sugar structures were generated using MAFFT v7.508 [45] to supply an initial multiple sequence alignment, followed by Aclust (section 4.1) for distance matrix embedding and clustering. The trees were visualized in iTOL [60].

4.7 Oligosaccharide backbone similarity score

A similarity score function was developed that quantifies the number of identical subunits at both donor and acceptor ends of oligosaccharides, specifically positions [..., -2, -1, +1, +2, ...] with respect to the bond formation site (Figure 4). The minimum non-zero similarity score between a pair of oligosaccharides is 2, requiring identity at both positions -1 and +1. Thereafter the comparison extends by one position in each positive (+2, +3, ...) and negative (-2, -3, ...) chain directions, adding one to the score for each additional identical match, but terminating at the first non-identity or possible re-use of a backbone position.

To facilitate comparison, oligosaccharides sequences are translated from IUPAC nomenclature into symbols that represent elements of backbone geometry, only considering monomer dimension and stereochemistry of acceptor and anomeric donor carbon atoms, and ignoring sidechains and chemical modification (Fig. 9). Briefly, the monomer dimension is represented by a single letter P, F or L depending on whether the monomer sugar is a pyranose, furanose or is linear, respectively. Stereochemistry of the acceptor and donor carbon atoms is represented by the index number of the carbon position within the ring/monomer, followed by a single letter U, D or N depending on whether the linked oxygen atom is U (up=above the monomer ring), D (down=below the monomer ring), or N (neither above or below the ring). The N symbol is assigned in cases of conformational flexibility such as with alditols or C6 linkages. At present, in scoring the similarity of two thus translated residues, the entirety of the translation strings must be identical to achieve a score of +1. Further details and limitations will be presented elsewhere (G.P. Gippert, manuscript in preparation).

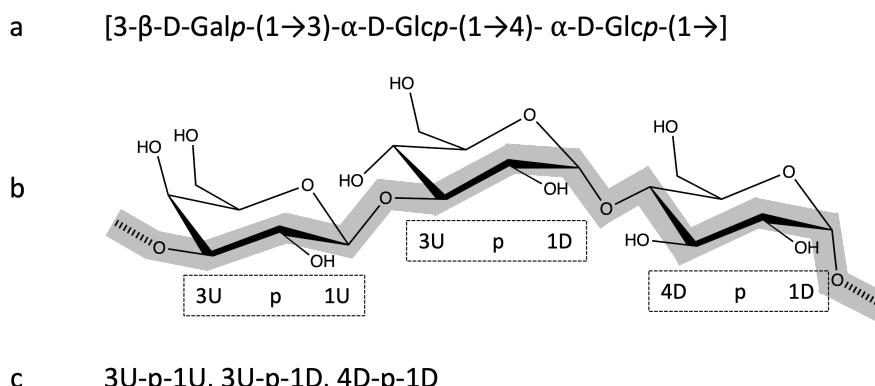


Figure 9: Oligosaccharide translation from IUPAC nomenclature to backbone (geometric) subunits for a trisaccharide consisting of one D-galactopyranose and two D-glucopyranose residues joined by intramolecular β 1 \rightarrow 3 and α 1 \rightarrow 4 bonds, respectively, and an intermolecular α 1 \rightarrow 3 bond formed in the polymerase reaction. (a) IUPAC nomenclature (b) Stereochemical projection highlighting backbone (thick grey line) and transfer bond (hatched line segments), and translated geometric subunits below (see text). (c) Completed translation.

474 4.8 Comparison of the families

475 Pairwise HHblits analyses [33] were performed for each of the new CAZy families. The HHblits scores were
476 visualized in a heatmap using Python Matplotlib [61].

477 AlphaFold2 [14] structures were generated of representative proteins from the families using the Colab-
478 Fold implementation [62] on our internal GPU cluster processed with the recommended settings. The best
479 ranked relaxed model was used. The protein structures were visualized in PyMOL [63] and pairwise structural
480 superimpositions were performed using the CEalign algorithm [64].

481 5 Data availability

482 Accessions to the seed sequences utilized in this work are given in Supplementary Table 1-2; the constantly
483 updated content of families GTxx1 - GTx17 is given in the online CAZy database at www.cazy.org.

484 6 Acknowledgements

485 This work was supported by the Novo Nordisk Foundation [grant number NNF20SA0067193]. Drs. Vincent
486 Lombard and Nicolas Terrapon are gratefully acknowledged for their assistance in incorporating our data into
487 the CAZy database. We also thank Dr. Philip Toukach for kindly providing a copy of the CSDB.

488 7 Author contributions

489 I.M. performed data acquisition, sequence analysis and interpretation; G.P.G. developed methodologies, super-
490 vised, analyzed and interpreted data; K.B. supervised, analyzed and interpreted data; C.J.H. performed custom
491 structure predictions; B.H. conceived the study, supervised and interpreted results. The manuscript was written
492 by I.M. and B.H. with help from all co-authors.

493 8 Competing interests

494 None

495 References

- 496 [1] Varki, A. *et al.* (eds.) *Essentials of Glycobiology* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor
497 (NY), 2022), 4th edn. URL <http://www.ncbi.nlm.nih.gov/books/NBK579918/>.
- 498 [2] Laine, R. A. A calculation of all possible oligosaccharide isomers both branched and linear yields 1.05×10^{12} structures for a reducing hexasaccharide: the Isomer Barrier to development of single-method
499 saccharide sequencing or synthesis systems. *Glycobiology* **4**, 759–767 (1994).
- 500 [3] Lapébie, P., Lombard, V., Drula, E., Terrapon, N. & Henrissat, B. Bacteroidetes use thousands of enzyme
501 combinations to break down glycans. *Nature Communications* **10**, 2043 (2019). URL <https://www.nature.com/articles/s41467-019-10068-5>.
- 502 [4] Lairson, L., Henrissat, B., Davies, G. & Withers, S. Glycosyltransferases: Structures, Functions, and
503 Mechanisms. *Annual Review of Biochemistry* **77**, 521–555 (2008). URL <https://www.annualreviews.org/doi/10.1146/annurev.biochem.76.061005.092322>.
- 504 [5] Drula, E. *et al.* The carbohydrate-active enzyme database: functions and literature. *Nucleic Acids Research*
505 **50**, D571–D577 (2022). URL <https://academic.oup.com/nar/article/50/D1/D571/6445960>.
- 506 [6] McDonald, A. G. & Tipton, K. F. Fifty-five years of enzyme classification: advances and difficulties. *The
507 FEBS journal* **281**, 583–592 (2014).
- 508 [7] Campbell, J. A., Davies, G. J., Bulone, V. & Henrissat, B. A classification of nucleotide-diphospho-sugar
509 glycosyltransferases based on amino acid sequence similarities. *The Biochemical Journal* **326** (Pt 3),
510 929–939 (1997).
- 511 [8] Coutinho, P. M., Deleury, E., Davies, G. J. & Henrissat, B. An Evolving Hierarchical Family Classification
512 for Glycosyltransferases. *Journal of Molecular Biology* **328**, 307–317 (2003). URL <https://linkinghub.elsevier.com/retrieve/pii/S0022283603003073>.

- 517 [9] Knauer, R. & Lehle, L. The oligosaccharyltransferase complex from yeast. *Biochimica et Biophysica Acta (BBA) - General Subjects* **1426**, 259–273 (1999). URL <https://linkinghub.elsevier.com/retrieve/pii/S0304416598001287>.
- 518
519
- 520 [10] Cho, H. Assembly of Bacterial Surface Glycopolymers as an Antibiotic Target. *Journal of Microbiology* (2023). URL <https://link.springer.com/10.1007/s12275-023-00032-w>.
- 521
522
- 523 [11] Sjodt, M. *et al.* Structure of the peptidoglycan polymerase RodA resolved by evolutionary coupling analysis. *Nature* **556**, 118–121 (2018). URL <http://www.nature.com/articles/nature25985>.
- 524
525
- 526 [12] Käshammer, L. *et al.* Cryo-EM structure of the bacterial divisome core complex and antibiotic target FtsWIQBL. *Nature Microbiology* **8**, 1149–1159 (2023). URL <https://www.nature.com/articles/s41564-023-01368-0>.
- 527
528
- 529 [13] Nygaard, R. *et al.* Structural basis of peptidoglycan synthesis by E. coli RodA-PBP2 complex. *Nature Communications* **14**, 5151 (2023). URL <https://www.nature.com/articles/s41467-023-40483-8>.
- 530
531
- 532 [14] Meeske, A. J. *et al.* SEDS proteins are a widespread family of bacterial cell wall polymerases. *Nature* **537**, 634–638 (2016). URL <http://www.nature.com/articles/nature19331>.
- 533
534
- 535 [15] Di Lorenzo, F. *et al.* A Journey from Structure to Function of Bacterial Lipopolysaccharides. *Chemical Reviews* **122**, 15767–15821 (2022). URL <https://pubs.acs.org/doi/10.1021/acs.chemrev.0c01321>.
- 536
537
- 538 [16] Islam, S. T. & Lam, J. S. Synthesis of bacterial polysaccharides via the Wzx/Wzy-dependent pathway. *Canadian Journal of Microbiology* **60**, 697–716 (2014). URL <http://www.nrcresearchpress.com/doi/10.1139/cjm-2014-0595>.
- 539
540
- 541 [17] Whitfield, C., Wear, S. S. & Sande, C. Assembly of Bacterial Capsular Polysaccharides and Exopolysaccharides. *Annual Review of Microbiology* **74**, 521–543 (2020). URL <https://www.annualreviews.org/doi/10.1146/annurev-micro-011420-075607>.
- 542
543
- 544 [18] Woodward, R. *et al.* In vitro bacterial polysaccharide biosynthesis: defining the functions of Wzy and Wzz. *Nature Chemical Biology* **6**, 418–423 (2010). URL <http://www.nature.com/articles/nchembio.351>.
- 545
546
- 547 [19] Bentley, S. D. *et al.* Genetic Analysis of the Capsular Biosynthetic Locus from All 90 Pneumococcal Serotypes. *PLoS Genetics* **2**, e31 (2006). URL <https://dx.plos.org/10.1371/journal.pgen.0020031>.
- 548
549
- 550 [20] Ruan, X., Loyola, D. E., Marolda, C. L., Perez-Donoso, J. M. & Valvano, M. A. The WaaL O-antigen lipopolysaccharide ligase has features in common with metal ion-independent inverting glycosyltransferases*. *Glycobiology* **22**, 288–299 (2012). URL <https://academic.oup.com/glycob/article-lookup/doi/10.1093/glycob/cwr150>.
- 551
552
- 553 [21] Ashraf, K. U. *et al.* Structural basis of lipopolysaccharide maturation by the O-antigen ligase. *Nature* **604**, 371–376 (2022). URL <https://www.nature.com/articles/s41586-022-04555-x>.
- 554
555
- 556 [22] Rai, A. K. & Mitchell, A. M. Enterobacterial Common Antigen: Synthesis and Function of an Enigmatic Molecule. *mBio* **11**, e01914–20 (2020). URL <https://journals.asm.org/doi/10.1128/mBio.01914-20>.
- 557
558
- 559 [23] Alexander, J. A. N. & Locher, K. P. Emerging structural insights into C-type glycosyltransferases. *Current Opinion in Structural Biology* **79**, 102547 (2023). URL <https://linkinghub.elsevier.com/retrieve/pii/S0959440X23000210>.
- 560
561
- 562 [24] Emami, K. *et al.* RodA as the missing glycosyltransferase in *Bacillus subtilis* and antibiotic discovery for the peptidoglycan polymerase pathway. *Nature Microbiology* **2**, 16253 (2017). URL <http://www.nature.com/articles/nmicrobiol2016253>.
- 563
564
- 565 [25] Maczuga, N., Tran, E. N. H., Qin, J. & Morona, R. Interdependence of *Shigella flexneri* O Antigen and Enterobacterial Common Antigen Biosynthetic Pathways. *Journal of Bacteriology* **204**, e00546–21 (2022). URL <https://journals.asm.org/doi/10.1128/jb.00546-21>.
- 566
567
- 568 [26] Lombard, V., Golaconda Ramulu, H., Drula, E., Coutinho, P. M. & Henrissat, B. The carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Research* **42**, D490–D495 (2014). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkt1178>.
- 569
570
- 571 [27] Servais, C. *et al.* Lipopolysaccharide biosynthesis and traffic in the envelope of the pathogen *Brucella abortus*. *Nature Communications* **14**, 911 (2023). URL <https://www.nature.com/articles/s41467-023-36442-y>.
- 572
573

- 566 [28] Iguchi, A. *et al.* A complete view of the genetic diversity of the *Escherichia coli* O-antigen biosynthesis gene cluster. *DNA Research* **22**, 101–107 (2015). URL <https://academic.oup.com/dnaresearch/article-lookup/doi/10.1093/dnares/dsu043>.
- 569 [29] Liu, B. *et al.* Structure and genetics of *Shigella* O antigens. *FEMS Microbiology Reviews* **32**, 627–653 (2008). URL <https://academic.oup.com/femsre/article-lookup/doi/10.1111/j.1574-6976.2008.00114.x>.
- 571 [30] Liu, B. *et al.* Structural diversity in *Salmonella* O antigens and its genetic basis. *FEMS Microbiology Reviews* **38**, 56–89 (2014). URL <https://academic.oup.com/femsre/article-lookup/doi/10.1111/1574-6976.12034>.
- 574 [31] Kenyon, J. J., Cunneen, M. M. & Reeves, P. R. Genetics and evolution of *Yersinia pseudotuberculosis* O-specific polysaccharides: a novel pattern of O-antigen diversity. *FEMS Microbiology Reviews* **41**, 200–217 (2017). URL <https://academic.oup.com/femsre/article/41/2/200/2996588>.
- 577 [32] Hu, D., Liu, B., Dijkshoorn, L., Wang, L. & Reeves, P. R. Diversity in the Major Polysaccharide Antigen of *Acinetobacter baumannii* Assessed by DNA Sequencing, and Development of a Molecular Serotyping Scheme. *PLoS ONE* **8**, e70329 (2013). URL <https://dx.plos.org/10.1371/journal.pone.0070329>.
- 580 [33] Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nature Methods* **9**, 173–175 (2012). URL <http://www.nature.com/articles/nmeth.1818>.
- 583 [34] Liu, B. *et al.* Structure and genetics of *Escherichia coli* O antigens. *FEMS Microbiology Reviews* **44**, 655–683 (2020). URL <https://academic.oup.com/femsre/article/44/6/655/5645236>.
- 585 [35] Li, C. *et al.* Structural, Biosynthetic, and Serological Cross-Reactive Elucidation of Capsular Polysaccharides from *Streptococcus pneumoniae* Serogroup 16. *Journal of Bacteriology* **201**, 13 (2019).
- 587 [36] Lin, F. L. *et al.* Identification of the common antigenic determinant shared by *Streptococcus pneumoniae* serotypes 33A, 35A, and 20 capsular polysaccharides. *Carbohydrate Research* **380**, 101–107 (2013). URL <https://linkinghub.elsevier.com/retrieve/pii/S000862151300284X>.
- 590 [37] Lin, F. L. *et al.* Structure elucidation of capsular polysaccharides from *Streptococcus pneumoniae* serotype 33C, 33D, and revised structure of serotype 33B. *Carbohydrate Research* **383**, 97–104 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621513003947>.
- 593 [38] Bush, C. A., Cisar, J. O. & Yang, J. Structures of Capsular Polysaccharide Serotypes 35F and 35C of *Streptococcus pneumoniae* Determined by Nuclear Magnetic Resonance and Their Relation to Other Cross-Reactive Serotypes. *Journal of Bacteriology* **197**, 2762–2769 (2015). URL <https://journals.asm.org/doi/10.1128/JB.00207-15>.
- 597 [39] Kondakova, A. N. *et al.* Reinvestigation of the O-antigens of *Yersinia pseudotuberculosis*: revision of the O2c and confirmation of the O3 antigen structures. *Carbohydrate Research* **343**, 2486–2488 (2008). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621508003443>.
- 600 [40] Henrissat, B. & Bairoch, A. Updating the sequence-based classification of glycosyl hydrolases. *The Biochemical Journal* **316** (Pt 2), 695–696 (1996).
- 602 [41] Paton, J. C. & Trappetti, C. *Streptococcus pneumoniae* Capsular Polysaccharide. *Microbiology Spectrum* **7**, 7.2.33 (2019). URL <https://journals.asm.org/doi/10.1128/microbiolspec.GPP3-0019-2018>.
- 604 [42] Doyle, L. *et al.* Mechanism and linkage specificities of the dual retaining β-Kdo glycosyltransferase modules of KpsC from bacterial capsule biosynthesis. *Journal of Biological Chemistry* **299**, 104609 (2023). URL <https://linkinghub.elsevier.com/retrieve/pii/S002192582300251X>.
- 607 [43] Krinos, C. M. *et al.* Extensive surface diversity of a commensal microorganism by multiple DNA inversions. *Nature* **414**, 555–558 (2001). URL <https://www.nature.com/articles/35107092>.
- 609 [44] Smith, T. & Waterman, M. Identification of common molecular subsequences. *Journal of Molecular Biology* **147**, 195–197 (1981). URL <https://linkinghub.elsevier.com/retrieve/pii/0022283681900875>.
- 611 [45] Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. *Molecular Biology and Evolution* **30**, 772–780 (2013). URL <https://academic.oup.com/mbe/article-lookup/doi/10.1093/molbev/mst010>.
- 614 [46] Sonnhammer, E. L. & Hollich, V. [No title found]. *BMC Bioinformatics* **6**, 108 (2005). URL <http://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-6-108>.

- 616 [47] Crippen, G. & Havel, T. *Distance Geometry and Molecular Conformation*, vol. 15 (Chemometrics Research
617 Studies Press Series, Research Studies Press, 1988).
- 618 [48] Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009). URL
619 <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-10-421>.
- 620 [49] Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching.
621 *Nucleic Acids Research* **39**, W29–W37 (2011). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr367>.
- 623 [50] Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide
624 sequences. *Bioinformatics* **22**, 1658–1659 (2006). URL <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btl158>.
- 626 [51] Huszczynski, S. M., Hao, Y., Lam, J. S. & Khursigara, C. M. Identification of the Pseudomonas aeruginosa
627 O17 and O15 O-Specific Antigen Biosynthesis Loci Reveals an ABC Transporter-Dependent Synthesis
628 Pathway and Mechanisms of Genetic Diversity. *Journal of Bacteriology* **202** (2020). URL <https://journals.asm.org/doi/10.1128/JB.00347-20>.
- 630 [52] Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function using networkx
631 (2008). URL <https://www.osti.gov/biblio/960616>.
- 632 [53] Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a
633 multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009). URL
634 <https://academic.oup.com/bioinformatics/article/25/9/1189/203460>.
- 635 [54] Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep protein annotation. *BMC
636 Bioinformatics* **20**, 473 (2019). URL <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3019-7>.
- 638 [55] Shannon, P. *et al.* Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction
639 Networks. *Genome Research* **13**, 2498–2504 (2003). URL <http://genome.cshlp.org/lookup/doi/10.1101/gr.1239303>.
- 641 [56] Petersen, B. O., Meier, S., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Determination of native
642 capsular polysaccharide structures of Streptococcus pneumoniae serotypes 39, 42, and 47F and comparison
643 to genetically or serologically related strains. *Carbohydrate Research* **395**, 38–46 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621514002560>.
- 645 [57] Petersen, B. O., Hindsgaul, O., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Structural elucidation
646 of the capsular polysaccharide from Streptococcus pneumoniae serotype 47A by NMR spectroscopy.
647 *Carbohydrate Research* **386**, 62–67 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621513004084>.
- 649 [58] Lam, J. S., Taylor, V. L., Islam, S. T., Hao, Y. & Kocíncová, D. Genetic and Functional Diversity of
650 Pseudomonas aeruginosa Lipopolysaccharide. *Frontiers in Microbiology* **2** (2011). URL <http://journal.frontiersin.org/article/10.3389/fmicb.2011.00118/abstract>.
- 652 [59] Toukach, P. V. & Egorova, K. S. Carbohydrate structure database merged from bacterial, archaeal, plant
653 and fungal parts. *Nucleic Acids Research* **44**, D1229–D1236 (2016). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv840>.
- 655 [60] Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display
656 and annotation. *Nucleic Acids Research* **49**, W293–W296 (2021). URL <https://academic.oup.com/nar/article/49/W1/W293/6246398>.
- 658 [61] Hunter, J. D. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* **9**, 90–95
659 (2007). Publisher: IEEE COMPUTER SOC.
- 660 [62] Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nature Methods* **19**, 679–682 (2022).
661 URL <https://www.nature.com/articles/s41592-022-01488-1>.
- 662 [63] Schrödinger, LLC. The PyMOL Molecular Graphics System, Version 1.8 (2015).
- 663 [64] Shindyalov, I. N. & Bourne, P. E. Protein structure alignment by incremental combinatorial extension
664 (CE) of the optimal path. *Protein Engineering* **11**, 739–747 (1998).