

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
6 bacterial polysaccharide polymerases (BP-Pol) are glycosyltransferases (GT) that build bacterial surface
7 polysaccharides. These integral membrane enzymes share the particularity of using diphospholipid-activated
8 sugars and were previously missing in the carbohydrate-active enzymes database (CAZy; www.cazy.org).
9 While the first three classes formed well-defined families of similar proteins, the sequences of BP-Pols were so
10 diverse that a single family could not be built. To address this, we developed a new clustering method where
11 a 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 ~~Photosynthesis has granted homotrophs access to vast amounts of carbohydrates which serve as abundant~~
19 ~~carbon sources for most heterotrophs.~~ Carbohydrate polymers (glycans) and glyco-conjugates ~~have thus become~~
20 ~~are~~ the most abundant biomolecules on Earth and adopt a wide range of functions including energy storage,
21 structure, signaling, and mediators of host-pathogen interactions [1]. Due to the stereochemical diversity of
22 monosaccharides and the many possible linkages they can engage into, glycans display an enormous structural
23 diversity [2, 3]. Yet, our knowledge on their assembly is far from complete, especially in comparison to the
24 enzymes catalyzing their ~~enzymatic~~ breakdown.

25 The transfer of sugar moieties to acceptor molecules such as proteins, lipids or other sugars, is ~~performed~~
26 ~~catalyzed~~ by enzymes called glycosyltransferases or GTs [4]. ~~GTs can be classified either by activity or by~~
27 ~~sequence similarity. The Enzyme Commission of the International Union of Biochemistry and Molecular Biology~~
28 ~~(IUBMB), has elaborated a classification system that integrates a description of the donor, acceptor and bond~~
29 ~~formed, summarized in the form of an EC number [5]. This activity-based classification, although enormously~~
30 ~~useful to avoid the proliferation of trivial names, has the limitation that it does not integrate the structural~~
31 ~~features of the enzymes nor can it easily accommodate enzymes that act on several substrates [5].~~

32 Campbell and colleagues (1997) proposed a sequence-based classification of GTs into 26 families, ~~which was~~
33 ~~subsequently expanded to 65 families in 2003 [6].~~ The number of sequence-based families has since continued to
34 grow based on the necessary presence of at least one experimentally characterized founding member to define
35 a family. ~~The constantly updated GT classification, and~~ is presented in the carbohydrate-active enzymes
36 database (CAZy; www.cazy.org) ~~along with similar family classifications of other carbohydrate-active enzymes~~
37 [7]. An ~~additional~~ advantage of the sequence-based classification is that it readily enables genome mining for the
38 presence of ~~new~~ family members. Today there are 116 GT families in the CAZy database and ~~this number will~~
39 ~~continue growing as novel glycosyltransferases are progressively discovered or as known GTs are incorporated~~
40 ~~in the database.~~ In contrast to the EC numbers [5], the sequence-based classification implicitly incorporates
41 the structural features of GTs including the conservation of the catalytic residues. ~~Structurally, there are two~~
42 ~~major folds for the nucleotide sugar dependent GTs, namely GT-A and GT-B, which both have Rossmann folds.~~
43 ~~By contrast, sugar phospholipid utilizing GTs are integral membrane proteins which have an overall GT-C fold~~
44 ~~with a number of transmembrane helices that varies from 8 to 13 [4].~~

45 It was recognized very early that sequence-based GT families group together enzymes that can utilize
46 different sugar donors and/or acceptors, illustrating how GTs can evolve to adopt novel substrates and form
47 novel products [8, 6]. Mechanistically, glycosyltransferases can be either retaining or inverting, based on the
48 relative stereochemistry of the anomeric carbon of the sugar donor and of the formed glycosidic bond [4]. ~~This~~
49 ~~With almost no exceptions, this~~ feature is conserved in previously defined sequence-based families, providing
50 predictive power to this classification, as the orientation of the glycosidic bond can be predicted ~~safely~~ even if
51 the precise transferred carbohydrate is not known.

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

60 ~~Bacteria synthesize various surface polysaccharides which confer them antigenic properties. Lipopolysaccharide~~
61 ~~(LPS) is a polysaccharide specific~~ The first of these four functional classes corresponds to the peptidoglycan
62 polymerases, SEDS (shape, elongation, division and sporulation) proteins. These proteins polymerize peptidoglycan
63 in complex with class B penicillin-binding proteins [10]. Several 3-D structures of SEDS proteins have been
64 determined, and they harbor 10 transmembrane helices and one long extracellular loop [11, ?, 12]. This loop
65 contains an Asp residue, which has been shown to be essential for SEDS function [11, 13].

66 The enzymes in the next two functional classes, bacterial polysaccharide polymerases (BP-Pol, also known
67 as Wzy) and O-antigen ligase (O-Lig, also known as WaaL) are involved in the synthesis of lipopolysaccharides
68 (LPS). LPS are polysaccharides on the membrane of Gram-negative bacteria, and consists of the serotype-specific
69 consist of the highly diverse O-antigen attached to the Lipid A-core oligosaccharide which is located in the
70 outer membrane [14]. On the other hand capsular polysaccharides (CPS also known as K-antigens) are
71 produced by both Gram-negative and Gram-positive bacteria [15]. The covalent anchoring of CPS is still
72 poorly understood, although it is found to be linked to peptidoglycan in some Gram-positives [15]. Bacteria from
73 the Enterobacterales order produce yet another type of surface polysaccharides referred to as the enterobacterial
74 common antigen (ECA), which consists of repeating units of N-acetylglucosamine, N-acetyl-D-mannosaminuronic
75 acid and 4-acetamido-4,6-dideoxy-D-galactose [16]. Most of these surface polysaccharides The structure of the
76 O-antigen determines the O-serotype of the bacteria. Most LPS structures are produced via the so-called
77 Wzx/Wzy-dependent pathway, which takes place on the plasma membrane (inner membrane in Gram-negatives)
78 [17, 15], for which the genes are located in a specific gene cluster [17]. In this pathway, sugar repeat units are
79 assembled on an undecaprenyl-diphosphate (BP-Pol catalyzes the polymerization of pre-assembled oligosaccharides
80 attached to Und-PP) anchor on the cytoplasmic side of the membrane and then flipped to the outside of
81 the membrane by the flippase Wzx. The repeat units are then polymerized by the bacterial polysaccharide
82 polymerases (Wzy, BP-Pols), by transferring the growing polymer to the incoming new repeat units [17, 18].
83 In the case of LPS, the Little is known about the activity of BP-Pols. Firstly, because they are difficult to
84 express heterologously, and to date, only one study has demonstrated the activity of O-Pol *in vitro* [18] and no
85 experimentally determined 3-D structure is available. Secondly, because the sequences of BP-Pols are highly
86 diverse with a sequence identity as low as 16% for different serotypes of the same species [17], it is difficult
87 to identify conserved residues. However, several studies have identified BP-Pols in the gene clusters of various
88 species, paving the way for analyzing BP-Pol sequences across a large range of taxonomic origin (see below).
89 These include some Gram-negative bacteria which also employ the Wzx/Wzy-dependent pathway to produce
90 capsular polysaccharides, including *Streptococcus pneumoniae* [19]. The third functional class, O-Lig catalyzes
91 the final step in the synthesis of LPS; the ligation of the newly synthesized polymer (O-antigen) is then ligated
92 onto Lipid A-core oligosaccharide by the O-antigen ligase (WaaL; O-Lig) [20]. ECA is produced via the same
93 pathway, but with another set of enzymes including the polymerase (WzyE)[20]. In order to distinguish these
94 polymerases from the serotype-specific polymerases, they are here referred to as ECA polymerases (ECA-Pols).
95 A structure of O-Lig in complex with Und-PP has been reported, which showed a fold with 12 transmembrane
96 helices and a long periplasmic loop containing several conserved residues; two Arg which bind to the phosphates
97 of Und-PP and a His which is proposed to activate the acceptor [21].

98 Several of The enzymes present in the fourth functional class, the enterobacterial common antigen polymerases
99 (ECA-Pol, also known as WzyE) are involved in the synthesis of enterobacterial common antigen (ECA). In
100 addition to the GTs from these pathways are missing from the CAZy database including ECA-Pols, BP-Pols O-antigen,
101 ECA is a specific polysaccharide that occurs on the cell surface in members of the Enterobacterales order. ECA
102 consists of repeating units of N-acetylglucosamine, and O-Ligs, as well as some peptidoglycan polymerases.
103 These enzymes share with CAZy family GT66 the particularity of catalyzing the transfer of oligosaccharides and
104 N-acetyl-D-mannosaminuronic acid and 4-acetamido-4-like GT66, their donor is also activated by a diphospholipid
105 (Und-PP). In an attempt to complete the sequence-based classification of GTs, we have performed a 6-dideoxy-D-galactose
106 [16]. ECA is also produced via the Wzy/Wzx-dependent pathway, where ECA-Pol performs the equivalent
107 reaction to the BP-Pols.

108 Structurally, the sugar-diphospholipid-utilizing GTs have an overall GT-C fold common to other integral
109 membrane GTs, which is different from the globular nucleotide-sugar-utilizing GTs: GT-A and GT-B [4]. GT-C
110 enzymes have a number of transmembrane helices that varies from 8 to 14 [4, 22]. Alexander and Locher recently
111 suggested two subgroups of GT-C glycosyltransferases, GT-CA and GT-CB [22], where O-Lig and SEDS make
112 up GT-CB [22]. As no structures have been published of ECA-Pol and BP-Pols, these have not been assigned

113 to a structural subgroup.

114 We have identified 17 new GT families covering a large number of the sugar-diphospholipid-utilizing GTs,
115 by detailed analysis of the primary sequence of peptidoglycan polymerases, polysaccharide polymerases and
116 O-antigen ligases to assign their sequences to CAZy families and SEDS proteins, ECA-Pols, BP-Pols and O-Ligs.
117 In addition, we examined how sequence diversity correlates with the diversity of the transferred oligosaccharides
118 and with the stereochemical outcome of the glycosyl transfer reaction. The analysis also revealed that the new
119 GT families organize in three clans across the functional classes suggestive of common ancestry. Despite of poor
120 sequence alignments we manage to identify conserved potentially critical amino acids common within the clans.
121

122 2 Results

123 2.1 Peptidoglycan Polymerases

124 The synthesis of peptidoglycan is primarily performed by class A penicillin binding proteins(PBPs), which
125 harbor a GT51 domain and a transpeptidase domain [23, 13]. However, it has been shown that peptidoglycan
126 polymerization is also performed by the proteins RodA [10] and FtsW [24], often called shape, elongation,
127 division and sporulation (SEDS) proteins. FtsW operates in complex with a transpeptidase that performs the
128 peptide cross linking [11]). For RodA and FtsWFor building the CAZy family of SEDS proteins, we used four
129 characterized proteins as seed sequences: the proteins with PDB IDs 6BAR [11], 8TJ3 [12] and 8BH1 [?], and
130 the protein with GenBank accession CAB15838.1 [25]. Family GTxx1 was created and initially populated by
131 using BLAST against GenBank, and subsequently by searching against GenBank with an HMM built from the
132 retrieved sequences. GTxx1 is a very large family currently counting over 57,200 GenBank members in the
133 CAZy database with a pairwise sequence identity of 19% over 221 residues for the most distant members.

134 The taxonomic distribution of family GTxx1 follows what was reported in [13], namely that this protein
135 family is present in all bacteria except for Mycoplasma. It is present in most but not all planctomycetes.

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

140 The three-dimensional structure of RodA from *Thermus thermophilus* has been determined and consists of
141 10 transmembrane helices with several large extracellular loops containing functionally important residues [11].
142 A large hydrophobic groove containing highly conserved residues is thought to be the lipid binding site. An
143 Asp residue has been shown to be essential for RodA function in both *T. thermophilus* and *B. subtilis* [11, 13].

144 Sequence-wise we found excellent sequence similarity between RodA and FtsW proteins from various sources
145 and they were easily grouped together in a single, very large family (GTxx1) currently counting over 57,200
146 members in the CAZy database and showing no significant sequence similarity to other GT families.

147 The taxonomic distribution of family GTxx1 follows what was reported in [13], namely that this protein
148 family is present in all bacteria except for Mycoplasma. It is present in most but not all planctomycetes.

149 2.2 Enterobacterial common antigen polymerases

150 The ECA-Pol which was studied in [26] was used as seed sequence for building the ECA-Pol family. Although the
151 CAZy database only lists Genbank GenBank entries [27], we decided to build our multiple sequence alignments
152 (MSAs) with sequences from the NCBI non-redundant database in order to capture more diversity. An ECA-Pol
153 sequence library was thus constructed from the seed sequence using BLAST against the non-redundant database
154 . The of the NCBI. The ECA-Pols display a high sequence conservation, consistent with the conservation of
155 acceptor, donor and product of the reaction. ECA-Pols were therefore assigned to a single new and homogeneous
156 CAZy family CAZy family, GTxx2. To date this new family contains over 4800 members. The repeat unit being
157 axially bound to Und-PP and axially linked in the final polymer, this reaction is retaining the configuration of
158 the anomeric carbon undergoing catalysis GenBank members with sequence identity greater than 38% over 414
159 residues, consistent with the conservation of acceptor, donor and product of the reaction.

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

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

165 2.3 O-antigen ligases

166 With the aim of including the O-Ligs in the CAZy database, we collected 37 O-Lig sequences (Supplementary
 167 Table 1) and constructed a sequence library from these seed sequences using BLAST against the NCBI non-
 168 redundant database. A phylogenetic tree was constructed with the sequence library using our in-house Aclust
 169 tool which revealed four distantly related clades (Supplementary Fig. 1). The O-Ligs were included into one
 170 new CAZy family, GTxx3 with >more than 16,700 members distributed in the four subfamilies.

171 The greater diversity of the GTxx3 O-antigen ligases O-Ligs compared to the GTxx1 peptidoglycan poly-
 172 merases and GTxx2 ECA-Pol appears in the form of the four divergent clades in the O-Lig phylogenetic tree
 173 (Supplementary Fig. 1). We hypothesize that this increased diversity originates from the extensive donor
 174 and moderate acceptor variability of O-Ligs [14]. Taxonomically, the GTxx3 O-Lig family is present in most
 175 bacteria, including both Gram-negatives and Gram-positives Gram-positive bacteria. The reaction performed
 176 by O-Ligs involves an inversion of the stereochemistry of the anomeric carbon since the sugar donor is axially
 177 bound to Und-PP and the reaction product is equatorially bound to Lipid A [20].

178 A recently discovered O-antigen ligase O-Lig, WadA, is bimodular with a GTxx3 domain appended to a
 179 globular glycosyltransferase domain of family GT25, which adds the last sugar to the oligosaccharide core [28].
 180 We have constructed a tree with representative WadA homologs from the GTxx3 family (Supplementary Fig.
 181 2) and observe that most of the sequences appended to a GT25 domain cluster together in one area, which form
 182 one clade in the tree, except for a few outliers. This suggests a coupled action of the GT25 and of the GTxx3 at
 183 least for the bimodular O-antigen ligases and maybe O-ligs and possibly for the entire family. The bimodular
 184 WadA ligase O-Lig is observed in five genera including Mesorhizobium and Brucella.

185 2.4 Other bacterial polysaccharide polymerases

186 We collected 365 bacterial The fourth functional subgroup of GT-CB are the BP-Pols. As previously mentioned,
 187 there is only one experimentally characterized BP-Pol sequences (also referred to as Wzy) from seven different
 188 studies that have reported [18], but several studies have identified BP-Pols from the polysaccharide gene clusters,
 189 and we decided to build our families based on these published reports. We thus collected 363 predicted BP-Pol
 190 sequences from sequences from seven studies for various species, both Gram-negatives and Gram-positives: Es-
 191 cherichia coli [29], Shigella boydii, Shigella dysenteriae, Shigella flexneri [30], Salmonella enterica [31], Yersinia
 192 pseudotuberculosis, Yersinia similis [32], Pseudomonas aeruginosa [17], Acinetobacter baumanii, Acinetobacter
 193 nosocomialis [33] and Streptococcus pneumoniae [19] (Supplementary Table 2).

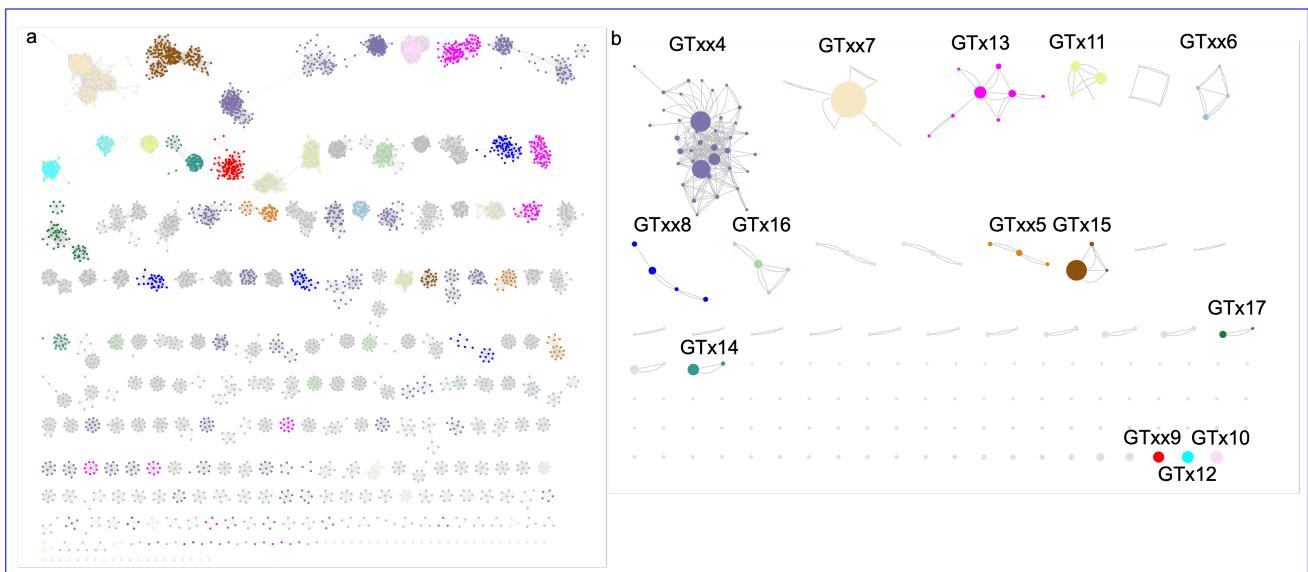


Figure 1: Clustering of BP-Pol sequences. a) SSN network with nodes representing proteins and edges representing pairwise alignment bit scores. b) HHblits network with nodes representing SSN clusters and edges representing HHblits scores. The resulting clusters are referred to as “superclusters”. There 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.

194 In contrast to ECA-Pols, the donors as well as the acceptors of BP-Pols are highly variable. Others have
 195 reported an exceptional sequence diversity of BP-Pols even within the same species [17]. We also found that

196 the sequences of BP-Pols are extremely diverse, and global alignments failed to reveal any conserved residue
197 due to both sequence diversity and to the difficulty in aligning proteins with multiple and variable numbers of
198 **hydrophobic-transmembrane** helices. It was therefore not possible to build a single family that could capture
199 the diversity of BP-Pols.

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

206 Instead, we used a combination of SSN and HMM comparisons: First, we used an SSN with a strict threshold
207 which would allow us to build good MSAs for the resulting clusters. This resulted in 204 clusters (Fig. 1a).
208 Next, we created an HMM profile of each SSN cluster and compared the HMMs by all-vs-all pairwise HHblits,
209 a program that aligns two HMMs and calculates a similarity score [34]. We then combined the SSN clusters
210 into **"superclusters"** "superclusters" in a network analysis based on the HHblits scores (Fig. 1b). **For this,**
211 **we used a score cut off of 160 in order to get a meaningful sequence and organismal diversity,** resulting in
212 28 **"superclusters"** "superclusters" of varying sizes and 86 singleton clusters. Interestingly, the BP-Pols **eluster**
213 **clustered** across taxonomy, and even BP-Pols from Gram-positive and Gram-negative bacteria **eluster** clustered
214 together. The 14 largest **"superclusters"** have been included as **superclusters define** new GT families in the
215 CAZy database (GTxx4-GTx17) with a number of members ranging from 159 to 5,979 at the time of submission.
216 Only **152 of the 365** **150 of the 363** original seeds are included in the new families. We thus expect that many
217 more BP-Pol families will be created in the future, as the amount and diversity of data increase.

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

221 **As a way of evaluating our families, we performed structural superimpositions of AlphaFold models of**
222 **distantly related members of each family. As an example, five distantly related members of GTxx4 are shown**
223 **in Supplementary Fig. 4. The sequence identity between these members is relatively low (between 21.4 and**
224 **24.3%). Yet, they still produce a meaningful superimposition, and notably, the conserved residues are oriented**
225 **very similarly.**

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. 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.

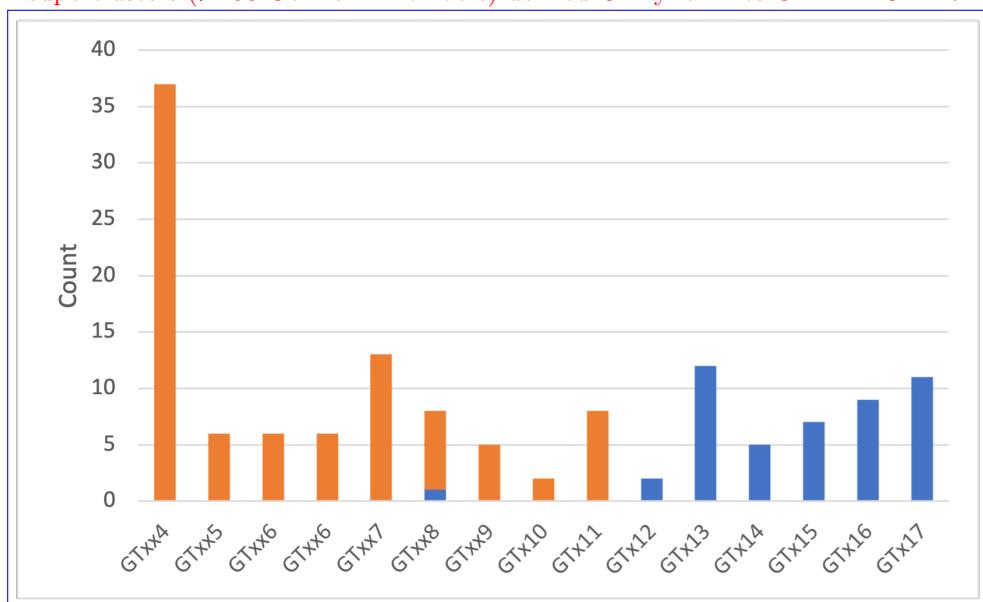


Figure 2: **Conservation of stereochemical outcome of the reaction catalyzed in the various BP-Pol families.** Equatorial bonds are shown in orange, implying an inverting mechanism. Axial bonds are shown in blue, implying a retaining mechanism.

226 **2.5 Analyzing the sugars transferred by polysaccharide polymerases**

227 **2.5 Analyzing the sugars transferred by bacterial polysaccharide polymerases**

228 There are two possible outcomes for the BP-Pol catalyzed polymerization reaction, either retaining or inverting
229 the α -configuration of the anomeric carbon of the carbohydrate carrying the Und-PP moiety. Examination of
230 the structure of the polysaccharides produced thus often reveals the stereochemistry of the bond formed by
231 the polymerases. In order to assess if the stereochemical mechanism is conserved in the families, we retrieved
232 the Next, we investigated how the BP-Pol families relate to the structures of the transferred sugar repeat
233 units from the Carbohydrate Structure Database (CSDB) [35]. As mentioned above, 152 of the original 365
234 original oligosaccharide repeat units. We retrieved the serotype-specific sugar structures, which were reported
235 in the review papers [36, 30, 31, 32, 17, 33, 19]. Additionally, nine sugar structures were included, which were
236 published after the review papers [37, 38, 39, 40, 41]. Out of the 150 BP-Pol seed sequences that were included
237 in the new families. Out of these 152 BP-Pols, 132 were matched CAZy families, we matched 131 with a
238 sugar structure. In these structures, the The repeat units are oligosaccharides with 3-7 monomers within the
239 backbone, often with branches. In several of the studies from which the BP-Pol sequences were retrieved most
240 of the cases, the bond which is formed by the polymerase has been identified [33, 32, 19, 31]. In cases where the
241 polymerase linkage was not clear from the literature, we identified it by comparing with similar sugar structures
242 from similar polymerases. in the review papers based on the other GTs in the gene cluster which assemble the
243 repeat units.

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

251 We found that the stereochemical outcome of BP-Pols appears well conserved within the new BP-Pol CAZy
252 families and varies from one family to another (Fig. 2). There are two apparent exceptions, however, where one
253 of the polymerase linkages has a different stereochemistry than the rest of the family. This is attributable to
254 when a wrong polysaccharide was assigned to the polysaccharide gene cluster comprising the polymerase (for
255 example if the bacteria produces several surface polysaccharides), or when there was is only one exception; in
256 family GTxx8, the polymerase linkages are all equatorial except for the O-antigen in *Pseudomonas aeruginosa*
257 O4, where it is axial. It is possible that there could an error in the chemical structure reported for the
258 polysaccharide, or when the linkage made by the polymerase was wrongly predicted. For example in family
259 GTxx4, one of the polymerase linkages is axial while the other 37 are equatorial or that the serotype designation
260 was incorrect or that the *P. aeruginosa* O4 polymerase constitutes an exception.

261 Next, we investigated whether there was a correlation between the structures of the transferred sugars
262 and the sequence similarity of the BP-Pols. We created phylogenetic trees of the BP-Pols in each family and
263 visualized them with the corresponding transferred repeat units. We observe that the sugars within each family
264 show similarity and this similarity appears to correlate with the structure of the tree (Fig. 23, Supplementary
265 Fig. 3). It seems likely that there is either an error in the chemical structure or that the bacteria contains
266 two different polymerases, one that catalyzes the formation of an equatorial bond and one that catalyzes the
267 formation of an axial bond4). The ends of the repeat units, ie. the subsite moieties immediately upstream (+1)
268 and downstream (-1) of the newly created bond (Fig. 4) seem to be most conserved whereas more variability
269 occurs in the middle part. We hypothesize that the +1 and -1 subsites are the moieties most important for
270 recognition by the active site of the BP-Pol.

271 Conservation of stereochemical outcome of the reaction catalyzed in the various BP-Pol families. The bond
272 formed by the polymerase was retrieved from literature or deduced by comparison to similar sugars from similar
273 polymerases. Equatorial bonds are shown in orange and axial bonds are shown in blue. We observe examples
274 of BP-Pols from distant taxonomic origin that cluster in the same CAZy family and have highly similar sugars.
275 For example, *Escherichia coli* O178 and *Streptococcus pneumoniae* 47A in GTxx7 transfer sugars with almost
276 identical backbones, suggestive of horizontal gene transfer. There is only a slight variance in the middle of
277 the repeat unit. This suggests that there is less constraints on the central part of the repeat unit than on the
278 extremities that define the donor and the acceptor.

279 To We next attempted to quantify the correlation between BP-Pol sequence and sugar structure, we analyzed
280 the oligosaccharide repeat units associated with each of the CAZy familiescarbohydrate structure. For this
281 purpose, we developed an original pairwise oligosaccharide similarity score.

282 In our scoring scheme, the similarity of two glycans is estimated by examining subsite moieties immediately
283 upstream and downstream of the newly created interosidic bondthe -1 and +1 subsites, as we hypothesize
284 expect that these are the moieties most fitting the active site of the polymerase BP-Pol (Fig. 34). The

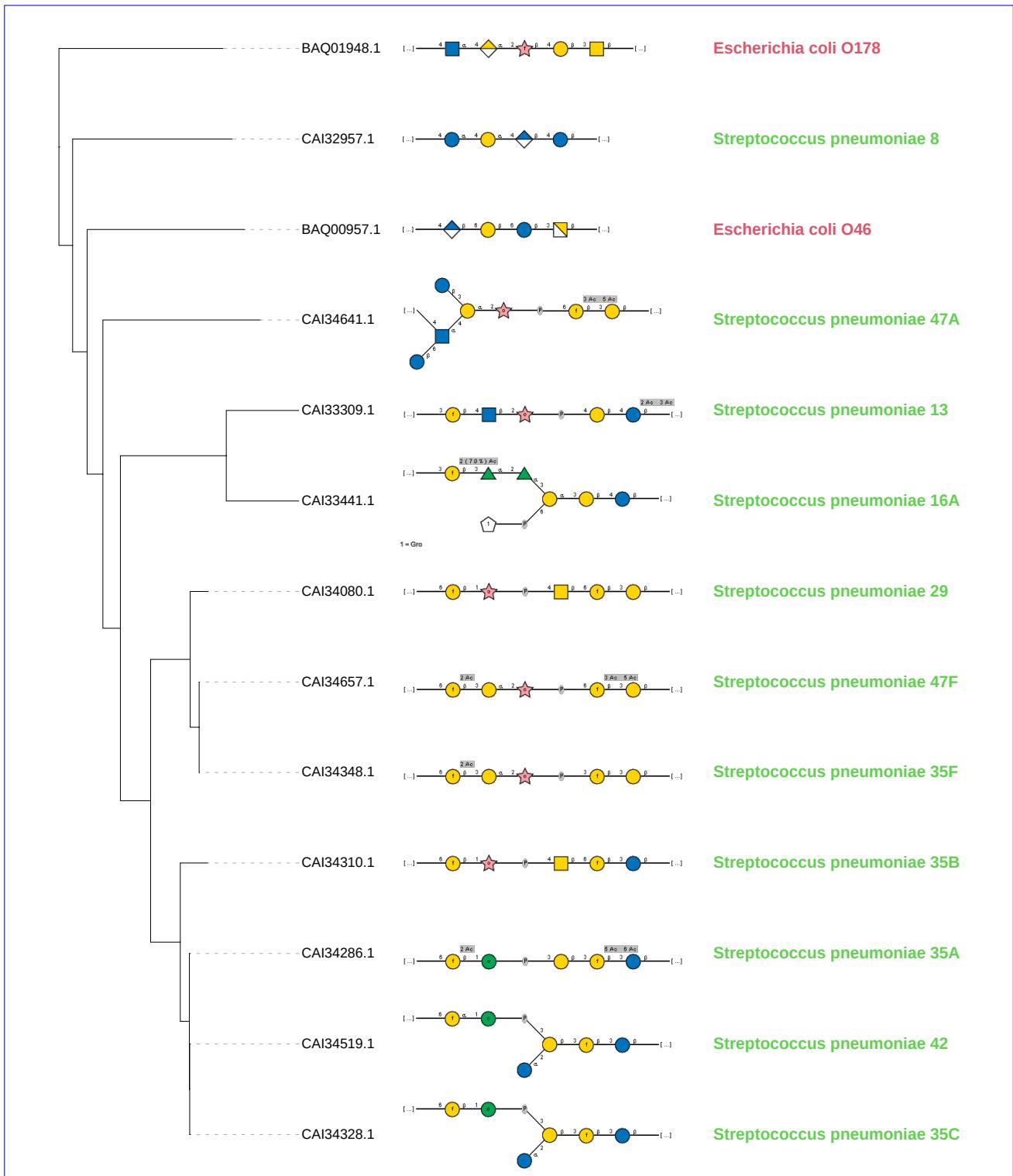


Figure 3: Similarity of transferred sugars by BP-Pols in GTxx7. The transferred repeat unit structures (in SNFG representation) are shown on a phylogenetic tree of BP-Pols in family GTxx7. There is an overall similarity between all the transferred sugars in the family and the similarity appears to correlate 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.

minimum match between two oligosaccharides corresponds to identical moieties at both subsites -1 and +1, which yields a score of 2. Thereafter, the score increases by one unit for each additional match at contiguous subsites, -2, -3, etc., and +2, +3, etc., up to a maximum value of 7 subsites found for the glycans encountered in this study (for details see Methods).

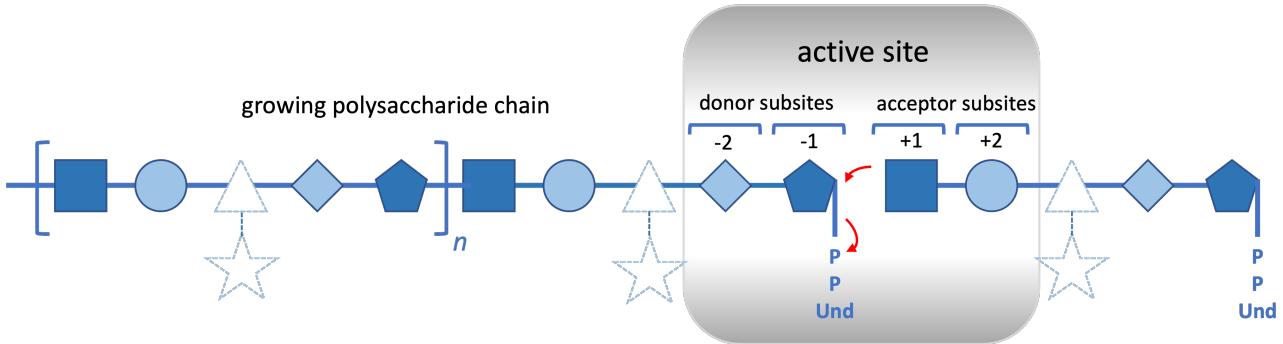


Figure 4: An idealized representation of a BP-Pol. The donor is the growing glycan chain activated by undecaprenyl pyrophosphate while the acceptor is a repeat unit monomer. The reaction is hypothesized to chiefly involve the sugar residues of the donor (subsites -2 and -1) and of the acceptor (subsites +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.

Using our glycan similarity scoring system, we found a correlation between sugar similarity and polymerase sequence similarity (Fig. 4)(Detailed: Supplementary Fig. 3)5), supported by a preponderance of similarity scores appearing close to the score matrix diagonal and within each individual family.

~~Notably, we observe examples of BP-Pols from distant taxonomic serotypes that cluster in the same CAZy family and have highly similar sugars. In Fig. 5, we show two examples of that. In GTxx7, two BP-Pols from very distant taxonomical origin (*Escherichia coli* and *Streptococcus pneumoniae*) transfer sugars with almost identical backbones. There is only a slight variance in the middle of the repeat unit. In GTx15, three BP-Pols from three different genera transfer glycans that, although consisting of different monosaccharides, have identical composition when translated into the backbone geometry description (Fig 5). These could be examples of horizontal gene transfer.~~

~~Examples of potential horizontal gene transfer. Top: Two BP-Pols from family GTxx7 from distant taxonomies transfer similar sugars. Bottom: Three different BP-Pols from different genera transfer similar sugars. The glycans are shown in SNFG representation and backbone geometry descriptors.~~

2.6 Comparison of families

Others have previously reported sequence and structural similarity between RodASEDS, O-Lig and some BP-Pols [12, 22, 21, 13]. In order to investigate the relatedness of the new CAZy families, we compared the family HMMs by all-vs-all HHblits analyses [34] (Fig. 6). Strikingly, we observe that the retaining BP-Pol families cluster together on the heatmap together along with the retaining ECA-Pols, while the inverting BP-Pols form two distinct groups, one of them containing the inverting O-Ligs. The background noise between some inverting and the retaining enzymes retaining enzymes is likely due to the general conservation of the successive transmembrane helices, which is altered in the GTxx4-GTxx5-GTxx6 subgroup due to their different architecture (vide infra); on the other hand, peptidoglycan polymerases see below). Peptidoglycan polymerases, GTxx1, segregate away from the other families.

Alexander and Locher recently suggested two subgroups of GT-C glycosyltransferases, GT-CA and GT-CB, based on the structural features of several of these families [22]. In the CAZy database, clans have been defined for the glycoside hydrolases (GHs), which group together CAZy families with distant sequence similarity, similar fold, similar catalytic machinery and stereochemical outcome [42]. In extension of the report of the GT-CB class by Alexander and Locher ([22])[22], and based on the above-mentioned similarities between the new CAZy families, we can now define three clans within GT-CB: GT-CB₁ consisting of inverting BP-Pol families and O-Lig, GT-CB₂ consisting of retaining BP-Pol families and ECA-Pol, and GT-CB₃ 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-CB₂ and GT-CB₃ to the structural subclass GT-CB of Alexander and Locher [22]. In addition, we also present in Table 1 the families of GT-C glycosyltransferases that have not yet been assigned to a structural class.

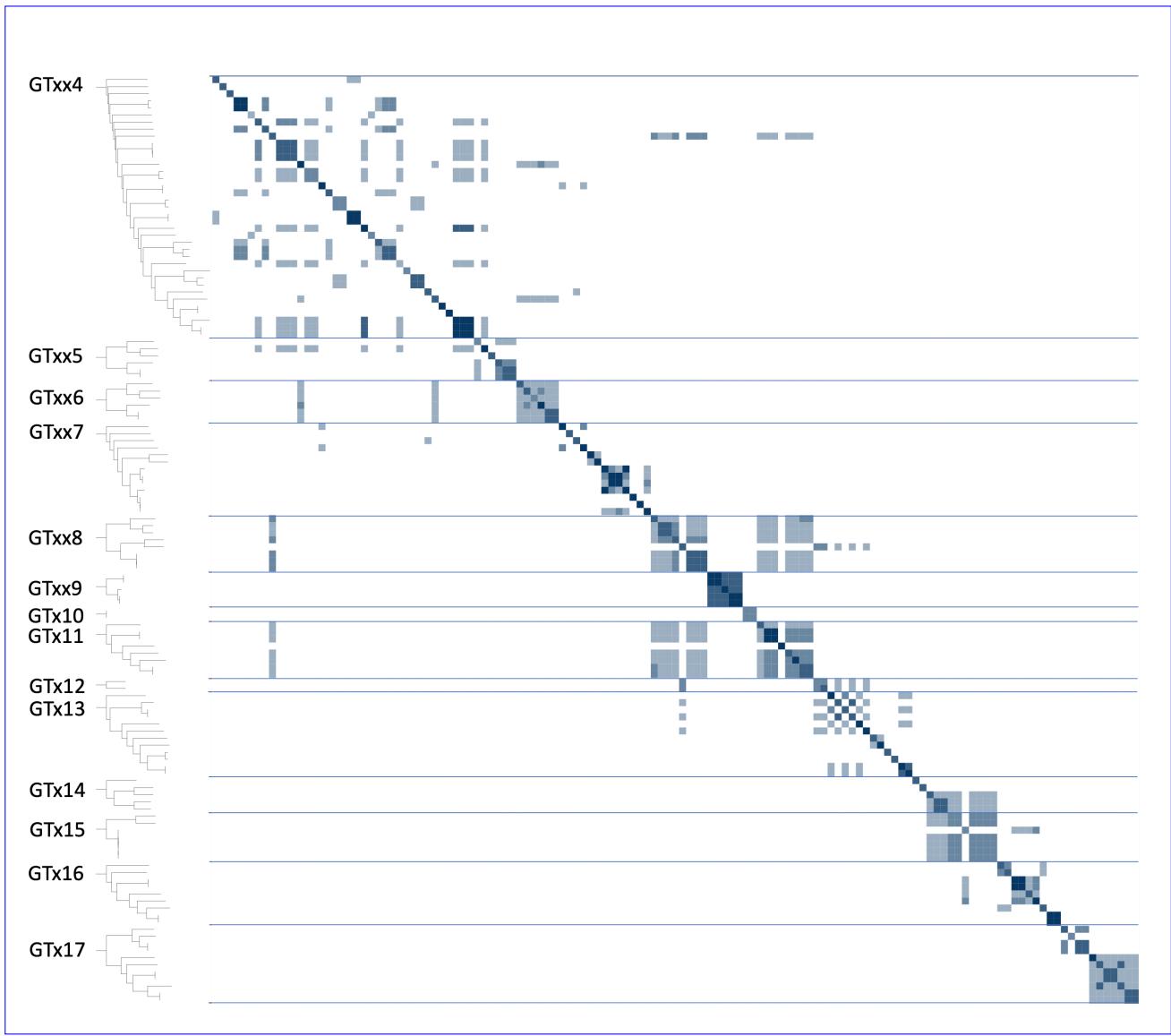


Figure 5: Glycan similarity of sugar repeat units polymerized by BP-Pols. All “seed” BP-Pols where the corresponding transferred oligosaccharide was known were included in the heatmap. A phylogenetic tree is shown for the polymerases in each CAZy family on the left. The glycan similarity scores are shown 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 [including both \$-2\$, \$+2\$ site for at least three additional sequential](#) positions). Blue lines separate the families.

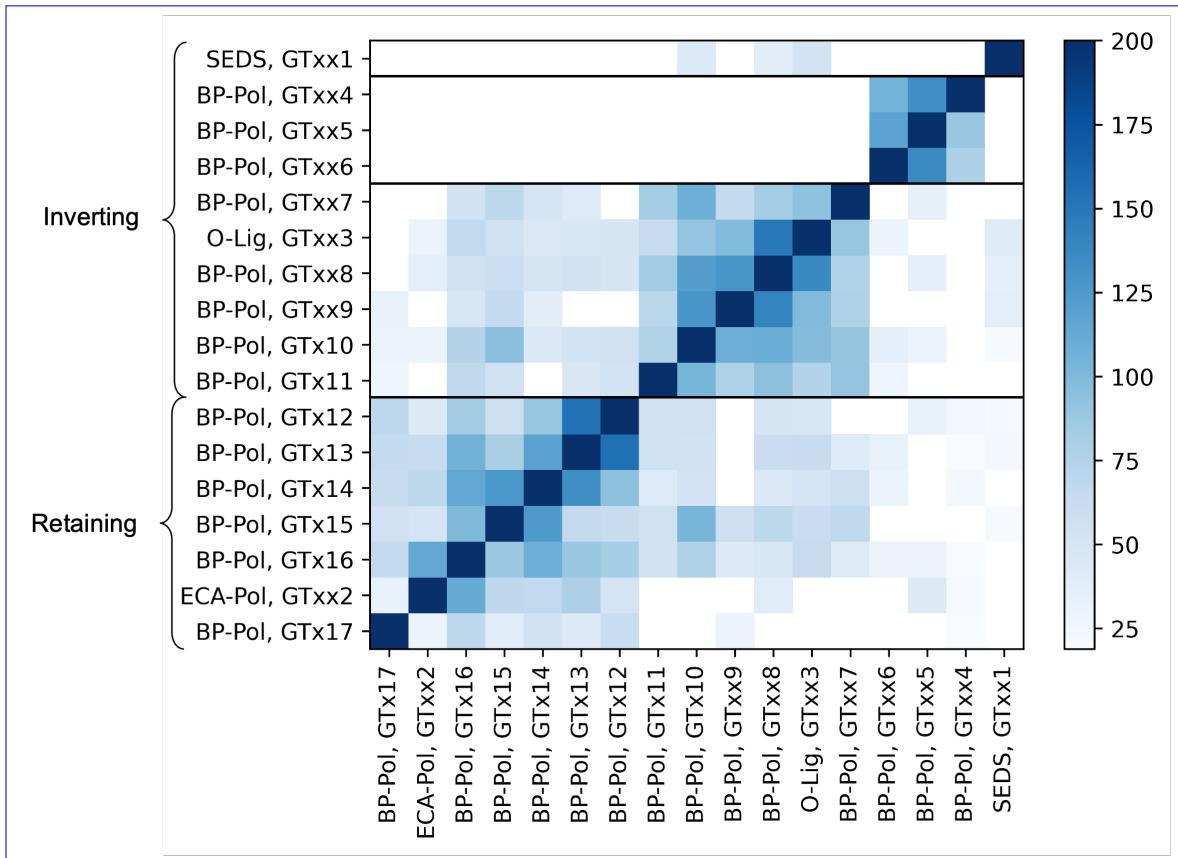


Figure 6: Heatmap of inter-family HHblits bit scores. The HHblits scores are shown on a color scale from white (low similarity score) to dark blue (high similarity score). The HHblits scores depend on the direction of the alignment, and therefore the heatmap is not symmetrical.

| 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.

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 on the outsidelocated on the outer face of the membrane. Enzymes of clans GT-CB1 and GT-CB2GT-CB1 and GT-CB2 have a long extracellular loop close to the C-terminus (Fig. 7). In stark contrast, families GTxx4, GTxx5 and GTxx6 of clan GT-CB3GT-CB3 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.

Most of 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-CB1 appear to have a similar pattern. All have 1-2 conserved Args, most of which are conserved in the familiesHHblits 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 the inverting Clan GT-CB1 have two conserved Arg residues and one conserved either Glu/Asp (in the BP-Pols) or His residue (in the O-Lig) (Fig. 7). In the pairwise HHblits alignments and structural superimpositions, the Glu/Asp /His residues align, suggesting that they could. 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 structural 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 O-Lig, the conserved His has been proposed to activate the acceptor, while the two Args are proposed to position the donor by binding to the phosphate groups [21]. We hypothesize that the Glu and Asp residues in the BP-Pols play the same role as the His in O-Lig.

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 are all in similar positions.

In the retaining clan GT-CB2, the pattern of conservation looks is different. Here, most of the families have 2-3 conserved Arg/Lys and one 1-2 conserved Tyr. As an example of the structural similarity in this clan, 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 the ECA-Pol family GTxx2 and family GTx16 is shown in Fig 8b. The structures again show each family shows that the conserved residues are oriented very similarly, despite the low overall similarity (RMSD 5.4 Å over 360 residues), but the conserved residues are oriented very similarly. This also shows that ECA-Pols display similarity to the BP-Pols of clan GT-CB2 (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-CB1. 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-CB1. We therefore hypothesize that this conserved Asp may play the role of activating the acceptor in clan GT-CB1 glycosyltransferases as the His in O-Lig [21].

3 Discussion

Here we have added 17 glycosyltransferase families (GTxx1 to GTx17) to the CAZy database bringing the total of covered families from 116 to 133. In the CAZy database, families are built by aggregating similar sequences around a biochemically characterized member. The known difficulties in the direct experimental characterization of integral membrane GTs render this constraint impractical. To circumvent this problem, but to remain connected to actual biochemistry, we decided to build our families around seed sequences for which knowledge of the glycosidic bond formed could be deduced from examination of the polysaccharide product from the literature. The list of these seed sequences is given in Supplementary Table 1-2 for families GTxx3 to GTx17. No seed sequence was needed for peptidoglycan polymerases (GTxx1) as the family is very tight around two structurally and functionally characterized members.

To our knowledge, this is the first time that BP-Pols from different species have been successfully clustered. Indeed, forming groups of BP-Pols has been very difficult before previously because of their extreme diversity even within strains of a single species [29], and, as a consequence, the knowledge on conserved and functional residues has been very limited. By combining BP-Pols from a wide range of taxonomical origins and expanding with the diversity from the NCBI non-redundant database current sequence diversity, we were able to form

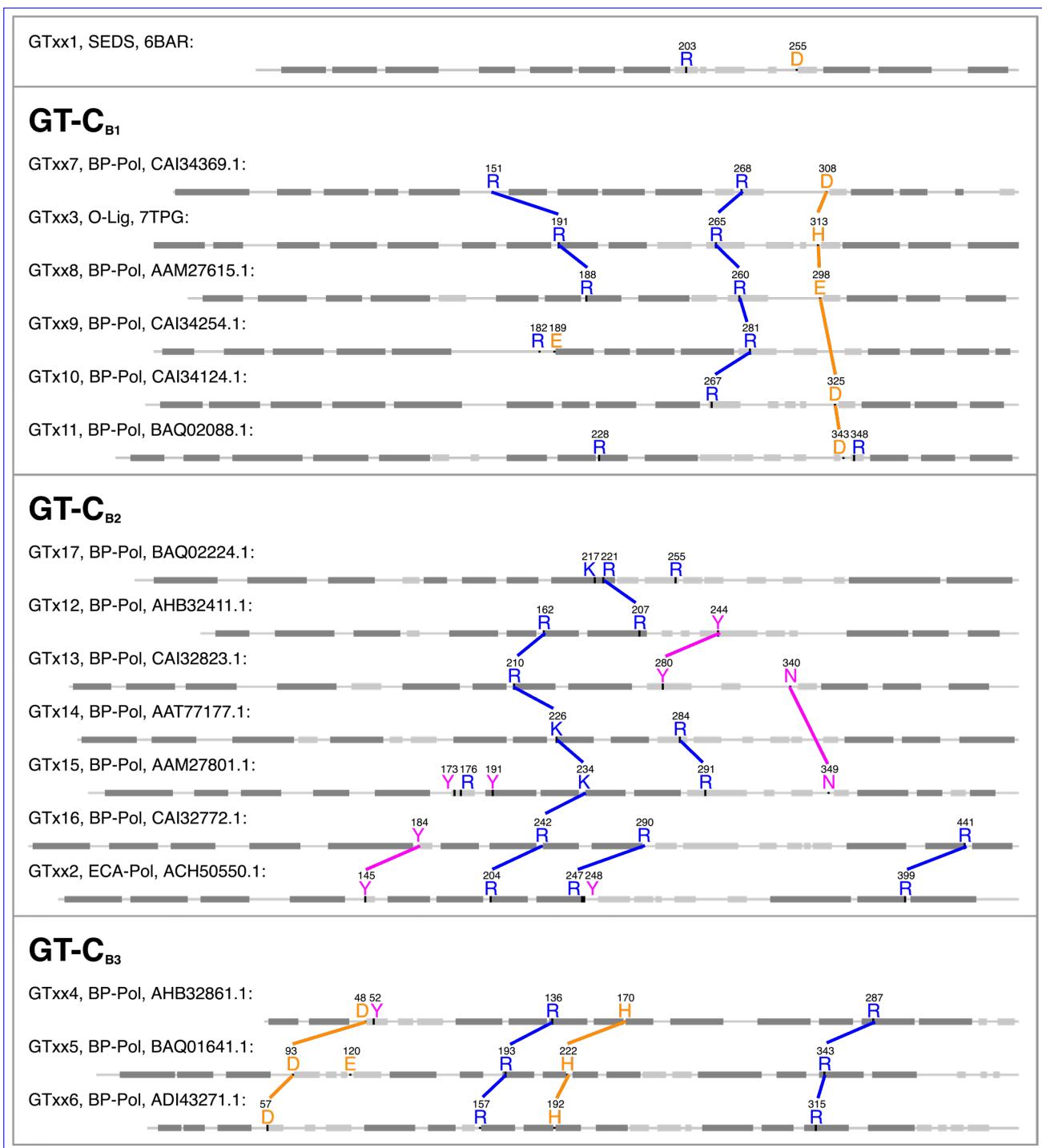


Figure 7: Comparison of equivalent conserved residues in the new GT families. The non-aliphatic conserved residues of each of the new CAZy families are shown on representative sequences of representative family members. Lines are shown between conserved residues that align in HHblits alignments and that 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. Lines are shown between residues that align in pairwise structural superimpositions. The secondary structure was retrieved 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 from experimental structures where available R in the family. Conserved aliphatic residues are not shown.

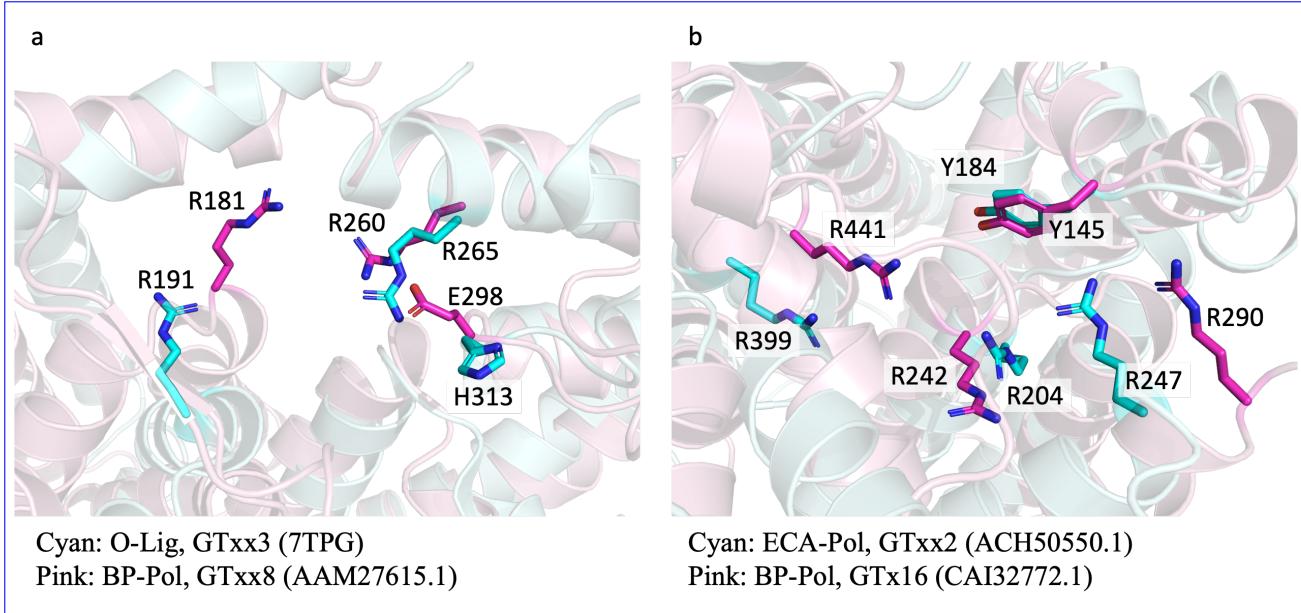


Figure 8: Structural superimposition of different families with conserved residues. a) O-Lig from GTxx3 (PDB 7TPG) and AlphaFold model of BP-Pol from GTxx8 (RMSD 5.3 Å over 192 residues). The conserved Glu in GTxx8 is aligning with the conserved His in GTxx3, which is proposed to activate the acceptor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

383 larger families of similar polymerases from widely different taxonomies, thereby revealing conserved residues
 384 that are most likely functionally important.

385 We observed that the O-Lig family (GTxx3) was present in many Gram-positive bacteria such as *Streptococcus pneumoniae*. ~~Gram-positive bacteria do not produce LPS, but instead capsular polysaccharides (CPS), which are linked to the peptidoglycan layer [43]~~ The covalent anchoring of CPS in Gram-negative bacteria is still poorly understood, although it is found to be linked to peptidoglycan in some Gram-positive bacteria [15, 43].
 386 Thus a hypothesis could be that the GTxx3 members in *S. pneumoniae* are responsible for the ligation of CPS to the peptidoglycan layer ~~in these bacteria~~.

387 Because families are more robust when built with enough sequence diversity, many clusters of O-antigen
 388 polymerases were judged too small to build meaningful CAZy families. Additional polymerase families are thus
 389 expected in the future with the accumulation of sequence data. For instance the small cluster that contains
 390 47% identical BP-Pols from *E. coli* (GenBank BAQ01516.1) and *A. baumanii* (GenBank AHB32586.1) only
 391 contains eight sequences and will remain unclassified until enough sequence diversity has accumulated. This
 392 arbitrary decision comes from the need to devise a classification that can withstand a massive increase in the
 393 number of sequences without the need to constantly revise the content of the families. ~~This new GT families
 394 based on O-antigen polymerases are poised to be formed when additional evidence becomes available.~~

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

399 It has been observed that for classical GT-A and GT-B fold glycosyltransferases, the catalytic mechanism is
 400 conserved within a family, but families with the same fold can have different mechanisms, possibly because the
 401 stereochemical outcome of the glycosyl transfer reaction is essentially dictated by the precise positioning and
 402 activation of the acceptor above ($\text{SNS}_{\text{N}}\text{N}_2$) or below ($\text{SNS}_{\text{N}}\text{N}_1$) the sugar ring of the donor [4]. Very occasionally,
 403 retaining glycosyltransferases have been shown to operate via a double displacement mechanism that involves
 404 Asp/Glu residues to form a glycosyl enzyme intermediate and to activate the acceptor that attacks this inter-
 405 mediate [44]. The families defined here display globally similar GT-C folds, and they also show conservation
 406 of the catalytic mechanism with about half of the families retaining and the other half inverting the anomeric
 407 configuration of the donor, suggesting that the outcome of the reaction catalyzed by GT-C glycosyltransferases
 408 is also dictated by the positioning of the acceptor with respect to the sugar plane of the acceptor. In turn this
 409 also suggests that retaining BP-Pols also operate by an $\text{SNS}_{\text{N}}\text{N}_1$ mechanism rather than by the formation of a
 410 glycosyl enzyme intermediate. This hypothesis is supported by the lack of invariant Asp or Glu residues which
 411 could be involved in the formation and subsequent breakdown of a glycosyl enzyme intermediate in the retain-
 412 ing families GTxx2 and GTx12-GTx17. Additionally, the $\text{SNS}_{\text{N}}\text{N}_1$ mechanism may provide protection against

417 the interception of a glycosyl enzyme intermediate by a water molecule resulting in an undesirable hydrolysis
418 reaction and termination of the polysaccharide elongation.

419 The wealth of structural data of GT-C glycosyltransferases now permits a deeper evaluation of the intrinsic
420 properties of this large class of enzymes. Alexander and Locher have recently evaluated the structural similarities
421 between GT-C fold glycosyltransferases and have divided them in two folding fold subclasses [22]. The GT
422 families that we describe here significantly expand the GT-C class in the CAZy database (www.cazy.org) and
423 allow to combine the structural classes with mechanistic information. Lairson *et al.* have proposed the
424 subdivision of GT-A and GT-B fold glycosyltransferases in clans that integrate the stereochemical outcome of
425 the reaction [4]. Here we also note the conservation of the stereochemistry in the families of BP-Pols and we
426 thus propose to group them into three clans which share the same fold, residual sequence conservation and
427 the same catalytic mechanism (Table 1). As more families of BP-Pols emerge, these three clans will likely
428 grow. Table 1 shows the three clans we defined here and how they relate to the structural classes defined by
429 Alexander and Locher. Of note are families GTxx4, GTxx5, and GTxx6 which do not bear any similarity, even
430 distant, with the GT families of the other two clans. These three families also stand out by the location in
431 the sequence of the long loop that harbors the catalytic site in the other GT-C families. In absence of relics of
432 sequence relatedness to the other families, GTxx4, GTxx5 and GTxx6 were assigned to clan GT-C_{B3}. With 10
433 transmembrane helices, it is tempting to suggest that this clan may belong to the folding fold subclass GT-C_B
434 of Alexander and Locher.

435 The analysis presented here shows that not only the stereochemistry of the glycosyl transfer is conserved
436 in the BP-Pol families, but our development of an original method to estimate glycan similarity also reveals
437 an unexpected a certain degree of structural similarity of the oligosaccharide repeat units, suggesting that the
438 latter constitutes a significant evolutionary constraint applying to the sequence and structure of BP-Pols. A
439 closer inspection of the oligosaccharide repeat units within the families further reveals that the carbohydrates
440 that appear the most constrained are the carbohydrates located (i) at the non-reducing end of the acceptor
441 and (ii) close to the undecaprenyl pyrophosphate of the donor, i.e. the residues closest to the reaction center
442 (Fig. 34). By contrast, residues away from the two extremities engaged in the polymerization reaction appear
443 more variable, and can tolerate insertions/deletions or the presence of flexible residues such as linear glycerol
444 or ribitol, with or without or the presence of a phosphodiester bond.

445 The version of the glycan similarity score presented here involves a direct was inspired in part by observed
446 structural similarities in different O-antigen repeat units assembled by very similar BP-Pols [17]. The repeat-unit
447 comparison involves a translation of glycan IUPAC nomenclature into terms representing to a reduced alphabet
448 of terms representing only backbone configuration, i.e., ignoring chemical modifications and sidechains. Furthermore,
449 a positive similarity score requires identical matches an entire identical match of all backbone elements
450 at both donor and acceptor positions (-1 and +1 sites in Fig. 34, respectively). Despite these simplifications,
451 the similarity score reveals, with exceptions, an overall greater intra- rather than inter-family oligosaccharide
452 similarity (Fig 5). These limitations will be addressed at a later stage (G.P. Gippert, in preparation).

453 We have next looked at the distribution of the new GT families in genomes, and particularly the families
454 of bacterial polysaccharide polymerasesBP-Pols. This uncovers broadly different schemes, with some bacteria
455 having only one polymerase (and therefore only able to produce a single polysaccharide) while others having
456 several, and sometimes more than 5, an observation in agreement with the report that Bacteroides fragilis
457 Bacteroides fragilis produces no less than 8 different polysaccharides from distinct genomic loci [45]. The multiplicity of polysaccharide biosynthesis loci in some genomes makes it sometimes difficult to assign a particular
458 polysaccharide structure to a particular biosynthesis operon. Although the families described here do not solve
459 all problems, their correlation with the stereochemical outcome of the glycosyl transfer reaction allows to resolve
460 some inconsistencies (vide supra).

461 As already shown in other occasions, the sequence-based classification of carbohydrate-active enzymes of
462 the CAZy database has predictive power. The case of the GT families described here supports this view as
463 the invariant residues in the families not only co-localize in the same area of the three-dimensional structures
464 (whether actual or AlphaFold-predicted), but also correspond to the residues found essential for function in
465 the families where this has been studied experimentally. The families described herein also show mechanistic
466 conservation and thus the stereochemistry of glycosyl transfer can be predicted. Finally, the observed similarity
467 in oligosaccharide repeat units that accompanies sequence similarity has also predictive power and paves the
468 way to the future possibility of in-silico in silico serotyping based on DNA sequence.
469

470 4 Methods

471 4.1 General methods used for building CAZy families

472 The sequence libraries for the different families were built from the seed sequences using “Blastp” from BLAST+
473 2.12.0+ [46] against the NCBI non-redundant database version 61. Redundancy reduction was performed using

474 CD-HIT 4.8.1 [47].
475 MSAs were generated with MAFFT v7.508 using the L-INS-i strategy (iterative refinement, using weighted
476 sum-of-pairs and consistency scores, of pairwise Needleman-Wunsh local alignments) [48]. HMMs were built
477 using the "hmmbuild" function from HMMER 3.3.2 [49]. The alignments were inspected in Jalview [50]. Finally,
478 the CAZy families were populated by a combination of the "hmmsearch" function from HMMER and Blastp
479 against Genbank.

480 4.1 Alignment-based Clustering (Aclust)

481 Phylogenetic trees were generated using an in-house tool called Aclust (G.P.Gippert, manuscript in preparation)
482 . Source code may be obtained via GitHub at <https://github.com/GarryGippert/Aclust>. Aclust employs a
483 hierarchical clustering algorithm comprising the following steps. (1) A distance matrix is computed from all-
484 vs-all pairwise local pairwise sequence alignments [51], or from a multiple sequence alignment provided by
485 MAFFT [48]. The distance calculation is based on a variation of Scoredist ([52]), however with distance values
486 normalized by [52] where distance values are normalized to the shorter pairwise sequence length rather than
487 to pairwise alignment length. (2) The distance matrix is embedded into orthogonal coordinates using metric
488 matrix distance geometry [53], and a(3) a bifurcating tree is computed using nearest-neighbor joining algorithm
489 is used to create an initial tree. (3and centroid averaging in the orthogonal coordinate space. The last centroid
490 created in this process is defined as the root node. (4) Beginning with the root node of the initial tree, each left
491 and right subtree constitutes disjoint subsets of the original sequence pool, which are embedded reembedded
492 and rejoined separately (i.e., step steps 2 and 3 repeated for each subset), and the process repeated recursively
493 — having the effect of gradually reducing deleterious effects on tree topology arising from long “long” distances
494 between unrelated proteins.

495 4.2 Building the peptidoglycan polymerase family (GTxx1)

496 The peptidoglycan polymerase family, GTxx1, was built by using Blastp from BLAST+ 2.12.0+ [46] with
497 the sequences of the characterized SEDS proteins (PDB 6BAR, 8TJ3 and 8BH1 and GenBank accession
498 CAB15838.1) against GenBank with a threshold of approximately 30% to retrieve the family members. Next,
499 an MSA was generated with MAFFT v7.508 using the L-INS-i strategy [48], and an HMM model was built
500 with hmmbuild of HMMER 3.3.2 [49]. The family was further populated using hmmsearch from HMMER 3.2.2
501 against GenBank.

502 4.3 Building the Enterobacterial common antigen polymerases family (GTxx2)

503 A sequence library of ECA-Pols was constructed by using “blastp” Blastp with the seed sequence (Genbank
504 GenBank accession AAC76800.1) against the NCBI non-redundant database as described in the section “General
505 methods used for building CAZy families”. All hits version 61 with an E-value smaller than threshold of 1e-
506 60 were selected. The . The hits were redundancy reduced using CD-HIT 4.8.1 [47] with a threshold of 99%.
507 The redundancy-reduced pool of ECA-Pol sequences was clustered using our in-house tool Aclust (see above),
508 and the tree showed one large clade and a few outliers. All the sequences in the large clade were used to build
509 the MSA for the family. The family was built and populated as described in the section “General methods used
510 for building CAZy families” an MSA using MAFFT v7.508 with the L-INS-i strategy [48]. An HMM was built
511 based on this MSA using hmmbuild of HMMER 3.3.2 [49]. The family GTxx3 was built in CAZy and populated
512 using Blastp against GenBank with an approximate threshold of 30% and hmmsearch against GenBank.

513 4.4 Building the O-antigen ligase family (GTxx3)

514 37 O-Lig sequences were selected from literature (Supplementary Table 1) and expanded using “blastp” Blastp
515 against the NCBI non-redundant database (see section “General methods used for building CAZy families”)
516 with an E-value cut-off of 1e-60, resulting in 13,431 hits. The blast hits were redundancy reduced. Redundancy
517 reduction was performed on the resulting sequence pool using CD-HIT with a threshold of 99%, resulting in
518 a pool of 1,402 sequences. A phylogenetic tree of the pool of O-Lig sequences was generated using Aclust
519 (see section 4.2 above), which showed deep clefts between main branches, and branches with sufficient internal
520 diversity (Supplementary Figure 2). Based on these results, four subfamilies were determined. An MSA was
521 built for the family as well as for the subfamilies , and the with MAFFT v7.508 using the L-INS-i strategy.
522 HMMs were built based on the MSAs using the hmmbuild of HMMER 3.3.2 [49]. The family was populated
523 as described in section “General methods used for building CAZy families” using Blastp against GenBank using
524 an approximate threshold of 30% identity with the seed sequences and using hmmsearch with the family and
525 subfamily HMMs.

526 4.5 Building the Bacterial polysaccharide polymerase families (GTxx4-GTx17)

527 365–363 BP-Pol sequences were selected from literature (from 2 phyla, 4 orders, 15 species; retrieved from review
528 papers on biosynthesis of O-antigens and capsular polysaccharides in different species: *Escherichia coli* [29],
529 *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri* [30], *Salmonella enterica* [31], *Yersinia pseudotuberculosis*,
530 *Yersinia similis* [32], *Pseudomonas aeruginosa* [17], *Acinetobacter baumannii*, *Acinetobacter nosocomialis* [33]
531 and *Streptococcus pneumoniae* [19] (complete list in Supplementary Table 2). The BP-Pols for *A. baumannii* O7
532 and O16 were omitted, because of uncertainty of their serotypes [33]. The BP-Pol from *P. aeruginosa* O15 was
533 also omitted, because it has been shown that this BP-Pol is inactivated and that the O-antigen is synthesized
534 via the ABC-dependent pathway rather than the Wzx/Wzy-dependent pathway [54].

535 The sequence library was expanded using “blastp” Blastp for each seed sequence against the NCBI non-
536 redundant database (46,644 hits). All hits with an E-value less than threshold of 1e-15 and a length between
537 320 and 600 residues were selected (29,372 hits). Redundancy reduction was performed using CD-HIT with a
538 threshold of 95% identity.

539 To build a find clusters of BP-Pol sequences that were large enough to create a CAZy family, we developed a
540 clustering method consisting of two steps. First, in order to make a sequence similarity network (SSN), all-vs-all
541 pairwise local alignments of the BP-Pol sequence pool, an all-vs-all pairwise local alignment was performed using
542 were performed using Blastp from BLAST+ 2.12.0+. The networks were visualized with Cytoscape [55]. A bit
543 score threshold of 110 was selected and the A series of networks were built using different bit score thresholds.
544 The members of the resulting SSN clusters were identified using NetworkX [56].

545 MSA and HMMs were and MSA of the members were built with MAFFT v7.508 using the L-INS-i strategy.
546 The MSAs were inspected using Jalview [50], and a bit score threshold of 110 was selected, as it was the lowest
547 score for which the SSN clusters had adequate sequence conservation (approximately 15 conserved residues).

548 HMMs were then built for each SSN cluster as described in section 4.1. The HMMs for each cluster using
549 hmmbuild of HMMER 3.3.2, and the HMMs were compared using HHblits 3.3.0 [57]. The HHblits network was
550 then visualized in Cytoscape [55] with an HHblits score threshold of 160. CAZy families A series of HHblits
551 networks were built using different HHblits score thresholds. Again, the members of the resulting “superclusters”
552 were identified using NetworkX and MSA of the members were built with MAFFT v7.508 using the L-INS-i
553 strategy. A bit score threshold of 160 was selected as it resulted in “superclusters” with adequate diversity for
554 building CAZy families (approximately 5 conserved residues). CAZy families were created for the 14 biggest
555 largest superclusters and populated with sequences present in Genbank as described in the section “General
556 methods used for building CAZy families”. GenBank by a combination of Blastp with the seed sequences and
557 hmmsearch. The networks were visualized with Cytoscape [55].

558 4.6 Analysis of sugar repeat-unit structures

559 A copy of the bacterial records in the CSDB database (<http://csdb.glycoscience.ru>) was provided by Philip
560 Toukach [35] and extracted into a listing of

561 4.6 Analysis of sugar repeat unit structures

562 In order to analyze the relation between BP-Pol seeds, linking NCBI protein accessions with CSDB entries
563 based on serotype. The repeat-unit structures were cross-checked with the literature. In cases where there were
564 several sugar structures for a serotype in CSDB and in the literature, we chose the candidate that was most
565 similar to sugar structures for related sequence and structure of the transferred repeat unit, we retrieved the
566 repeat unit structures for the serotypes for the BP-Pols that were included in the new CAZy families. The
567 repeat unit structures were retrieved from the same review papers from which we got the BP-Pol sequences
568 [33, 19, 32, 31, 30, 17], except for the sugars for *E. coli*, where the sugar structures have been reported elsewhere
569 [36]. Nine additional repeat unit structures were included for *S. pneumoniae*, which were published after the
570 review paper; serotypes 16A [37], 33A [38], 33C and 33D [39], 35C and 35F [40], 42 and 47F [58] and 47A
571 [59]. For *Y. pseudotuberculosis* O3 and *S. pneumoniae* 33B, we used the revised structure from [41] and [39]
572 respectively. *Pseudomonas aeruginosa* O2 and O16 contain two BP-Pol genes; one BP-Pol localized in the
573 O-antigen biosynthesis cluster, which polymerizes the sugar repeat units with an α bond and one BP-Pol
574 localized outside the biosynthesis cluster which polymerizes the repeat units with a β bond [60]. Since the
575 BP-Pol reported in [17] are the BP-Pols from the O-antigen cluster, we report the sugar structure with the α
576 bond.

577 It is often, but not always, known which bond of the polysaccharide is created by the polymerase. The
578 sugar structures in CSDB are thus shown as the repeat units acted upon by BP-Pol, i.e., the bond made. The
579 linkages formed by the polymerase is the bond between the rightmost monosaccharide (the -1 site position, see
580 also Fig 3.) and the leftmost monosaccharide (the +1 site position). However, there are has been determined
581 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 repeat unit is provided in another “phase”, ie. polymerase linkage has not been determined unambiguously are *E. coli* O166, O78, the bond predicted to be catalyzed by BP-Pol is positioned internally within the linear repeat unit rather than at an end. We cross-checked the “phases” with the literature, and in cases where this was not provided in the literature, we compared them to other sugar structures from related BP-Pol sequences 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 [33]. *S. pneumoniae* 47A has WcjG as the initial transferase, which transfers Galp or Galf [33]. 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 could predict the polymerase bond and rearrange the sugar structures manually to show the putative correct phase. proposed the equivalent polymerase linkage.

The CSDB database (<http://csdb.glycoscience.ru>) [35] was used to retrieve literature, SNFG image representations of the carbohydrates were generated at the CSDB website.

and linear sugar strings of the repeat unit structures. Phylogenetic trees for only seed sequences in each of the newly-created BP-Pol families with sugar structures were generated using MAFFT v7.508 [48] to supply an initial multiple sequence alignment, followed by Aclust (section 4.24.1) for distance matrix embedding and clustering. Seed sequences are those where the sugar repeat unit structure is known. The trees were visualized with the corresponding sugar structures in iTOL [61].

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 the scoring, we have chosen to first translate oligosaccharides comparison, oligosaccharides sequences are translated from IUPAC nomenclature into a set of simplified geometric subunits that represent only symbols that represent elements of backbone geometry, only considering monomer dimension and stereochemistry of acceptor and anomeric donor carbon atoms, thus focusing entirely on the glycan backbone 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 linear/openis 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), this latter category. The N symbol is assigned in the cases of extensive cases of conformational flexibility such as with alditols or C6 linkages. Chemical modifications, side chains, and the configuration of non-linking carbons are ignored. Further details, limitations and extensions 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).

4.8 Comparison of the families

Pairwise HHblits analyses [34] were performed for each of the new CAZy families. The HHblits scores were visualized in a heatmap using Python Matplotlib [62].

AlphaFold2 [13] structures were generated of representative proteins from the families using the ColabFold implementation [63] on our internal GPU cluster processed with the recommended settings. The best ranked relaxed model was used[63]. The protein structures were visualized in PyMOL [64] and pairwise structural superimpositions were performed using the CEalign algorithm [65].

4.9 AlphaFold structures

The included AlphaFold2 predicted structures were generated using the ColabFold implementation on our internal GPU cluster processed with the recommended settings using the best ranked relaxed model [63]. The protein structures were visualized in PyMOL [64] and structural superimpositions were performed using the CEalign algorithm [65].

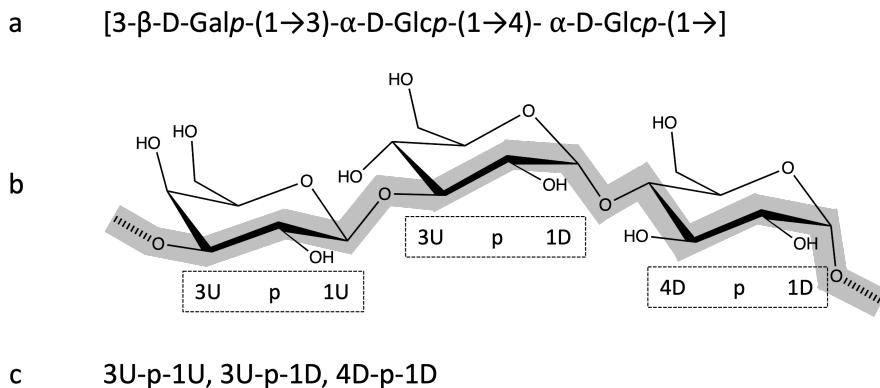


Figure 9: Oligosaccharide translation from IUPAC nomenclature to backbone (geometric) subunits for a trisaccharide consisting of one D-galactopyranose and two D-glucopyranose residues joined by intramolecular β 1 \rightarrow 3 and α 1 \rightarrow 4 bonds, respectively, and an intermolecular α 1 \rightarrow 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.

5 Data availability

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

6 Acknowledgements

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

7 Author contributions

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

8 Competing interests

None

References

- [1] Varki, A. et al. (eds.) *Essentials of Glycobiology* (Cold Spring Harbor Press, Cold Spring Harbor (NY), 2022), 4th edn. URL <http://www.ncbi.nlm.nih.gov/books/NBK579918/>.
- [2] Laine, R. A. A calculation of all possible oligosaccharide isomers both branched and linear yields 1.05 x 10(12) structures for a reducing hexasaccharide: the Isomer Barrier to development of single-method saccharide sequencing or synthesis systems. *Glycobiology* **4**, 759–767 (1994).
- [3] Lapébie, P., Lombard, V., Drula, E., Terrapon, N. & Henrissat, B. Bacteroidetes use thousands of enzyme combinations to break down glycans. *Nature Communications* **10**, 2043 (2019). URL <https://www.nature.com/articles/s41467-019-10068-5>.
- [4] Lairson, L., Henrissat, B., Davies, G. & Withers, S. Glycosyltransferases: Structures, Functions, and Mechanisms. *Annual Review of Biochemistry* **77**, 521–555 (2008). URL <https://www.annualreviews.org/doi/10.1146/annurev.biochem.76.061005.092322>.

- 663 [5] McDonald, A. G. & Tipton, K. F. Fifty-five years of enzyme classification: advances and difficulties. *The*
664 *FEBS journal* **281**, 583–592 (2014).
- 665 [6] Coutinho, P. M., Deleury, E., Davies, G. J. & Henrissat, B. An Evolving Hierarchical Family Classification
666 for Glycosyltransferases. *Journal of Molecular Biology* **328**, 307–317 (2003). URL <https://linkinghub.elsevier.com/retrieve/pii/S0022283603003073>.
- 668 [7] Drula, E. *et al.* The carbohydrate-active enzyme database: functions and literature. *Nucleic Acids Research*
669 **50**, D571–D577 (2022). URL <https://academic.oup.com/nar/article/50/D1/D571/6445960>.
- 670 [8] Campbell, J. A., Davies, G. J., Bulone, V. & Henrissat, B. A classification of nucleotide-diphospho-sugar
671 glycosyltransferases based on amino acid sequence similarities. *The Biochemical Journal* **326** (Pt 3),
672 929–939 (1997).
- 673 [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>.
- 676 [10] Cho, H. Assembly of Bacterial Surface Glycopolymers as an Antibiotic Target. *Journal of Microbiology*
677 (2023). URL <https://link.springer.com/10.1007/s12275-023-00032-w>.
- 678 [11] Sjodt, M. *et al.* Structure of the peptidoglycan polymerase RodA resolved by evolutionary coupling analysis.
679 *Nature* **556**, 118–121 (2018). URL <http://www.nature.com/articles/nature25985>.
- 680 [12] 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>.
- 682 [13] Meeske, A. J. *et al.* SEDS proteins are a widespread family of bacterial cell wall polymerases. *Nature* **537**,
683 634–638 (2016). URL <http://www.nature.com/articles/nature19331>.
- 684 [14] 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>.
- 686 [15] Whitfield, C., Wear, S. S. & Sande, C. Assembly of Bacterial Capsular Polysaccharides and Exopolysac-
687 charides. *Annual Review of Microbiology* **74**, 521–543 (2020). URL <https://www.annualreviews.org/doi/10.1146/annurev-micro-011420-075607>.
- 689 [16] Rai, A. K. & Mitchell, A. M. Enterobacterial Common Antigen: Synthesis and Function of an Enigmatic
690 Molecule. *mBio* **11**, e01914–20 (2020). URL <https://journals.asm.org/doi/10.1128/mBio.01914-20>.
- 691 [17] Islam, S. T. & Lam, J. S. Synthesis of bacterial polysaccharides via the Wzx/Wzy-dependent pathway.
692 *Canadian Journal of Microbiology* **60**, 697–716 (2014). URL <http://www.nrcresearchpress.com/doi/10.1139/cjm-2014-0595>.
- 694 [18] Woodward, R. *et al.* In vitro bacterial polysaccharide biosynthesis: defining the functions of Wzy and Wzz.
695 *Nature Chemical Biology* **6**, 418–423 (2010). URL <http://www.nature.com/articles/nchembio.351>.
- 696 [19] Bentley, S. D. *et al.* Genetic Analysis of the Capsular Biosynthetic Locus from All 90 Pneumococcal
697 Serotypes. *PLoS Genetics* **2**, e31 (2006). URL <https://dx.plos.org/10.1371/journal.pgen.0020031>.
- 698 [20] Ruan, X., Loyola, D. E., Marolda, C. L., Perez-Donoso, J. M. & Valvano, M. A. The WaaL O-antigen
699 lipopolysaccharide ligase has features in common with metal ion-independent inverting glycosyltrans-
700 ferases*. *Glycobiology* **22**, 288–299 (2012). URL <https://academic.oup.com/glycob/article-lookup/doi/10.1093/glycob/cwr150>.
- 702 [21] Ashraf, K. U. *et al.* Structural basis of lipopolysaccharide maturation by the O-antigen ligase. *Nature* **604**,
703 371–376 (2022). URL <https://www.nature.com/articles/s41586-022-04555-x>.
- 704 [22] Alexander, J. A. N. & Locher, K. P. Emerging structural insights into C-type glycosyltransferases. *Current
705 Opinion in Structural Biology* **79**, 102547 (2023). URL <https://linkinghub.elsevier.com/retrieve/pii/S0959440X23000210>.
- 707 [23] Goffin, C. & Ghysen, J.-M. Multimodular Penicillin-Binding Proteins: An Enigmatic Family of Or-
708 thologs and Paralogs. *Microbiology and Molecular Biology Reviews* **62**, 1079–1093 (1998). URL <https://journals.asm.org/doi/10.1128/MMBR.62.4.1079-1093.1998>.

- 710 [24] Taguchi, A. *et al.* FtsW is a peptidoglycan polymerase that is functional only in complex with its cog-
711 nate penicillin-binding protein. *Nature Microbiology* **4**, 587–594 (2019). URL <https://www.nature.com/articles/s41564-018-0345-x>.
- 713 [25] Emami, K. *et al.* RodA as the missing glycosyltransferase in *Bacillus subtilis* and antibiotic discovery for
714 the peptidoglycan polymerase pathway. *Nature Microbiology* **2**, 16253 (2017). URL <http://www.nature.com/articles/nmicrobiol2016253>.
- 716 [26] Maczuga, N., Tran, E. N. H., Qin, J. & Morona, R. Interdependence of *Shigella flexneri* O Antigen and
717 Enterobacterial Common Antigen Biosynthetic Pathways. *Journal of Bacteriology* **204**, e00546–21 (2022).
718 URL <https://journals.asm.org/doi/10.1128/jb.00546-21>.
- 719 [27] Lombard, V., Golaconda Ramulu, H., Drula, E., Coutinho, P. M. & Henrissat, B. The carbohydrate-
720 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>.
- 722 [28] 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>.
- 725 [29] 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>.
- 728 [30] Liu, B. *et al.* Structure and genetics of *Shigella* O antigens. *FEMS Microbiology Reviews* **32**, 627–653 (2008).
729 URL <https://academic.oup.com/femsre/article-lookup/doi/10.1111/j.1574-6976.2008.00114.x>.
- 730 [31] Liu, B. *et al.* Structural diversity in *Salmonella* O antigens and its genetic basis. *FEMS Microbiology*
731 *Reviews* **38**, 56–89 (2014). URL <https://academic.oup.com/femsre/article-lookup/doi/10.1111/1574-6976.12034>.
- 733 [32] Kenyon, J. J., Cunneen, M. M. & Reeves, P. R. Genetics and evolution of *Yersinia pseudotuberculosis* O-
734 specific polysaccharides: a novel pattern of O-antigen diversity. *FEMS Microbiology* **41**, 200–217
735 (2017). URL <https://academic.oup.com/femsre/article/41/2/200/2996588>.
- 736 [33] Hu, D., Liu, B., Dijkshoorn, L., Wang, L. & Reeves, P. R. Diversity in the Major Polysaccharide Antigen
737 of *Acinetobacter baumannii* Assessed by DNA Sequencing, and Development of a Molecular Serotyping
738 Scheme. *PLoS ONE* **8**, e70329 (2013). URL <https://dx.plos.org/10.1371/journal.pone.0070329>.
- 739 [34] Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence
740 searching by HMM-HMM alignment. *Nature Methods* **9**, 173–175 (2012). URL <http://www.nature.com/articles/nmeth.1818>.
- 742 [35] Toukach, P. V. & Egorova, K. S. Carbohydrate structure database merged from bacterial, archaeal, plant
743 and fungal parts. *Nucleic Acids Research* **44**, D1229–D1236 (2016). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv840>.
- 745 [36] Liu, B. *et al.* Structure and genetics of *Escherichia coli* O antigens. *FEMS Microbiology* **44**,
746 655–683 (2020). URL <https://academic.oup.com/femsre/article/44/6/655/5645236>.
- 747 [37] Li, C. *et al.* Structural, Biosynthetic, and Serological Cross-Reactive Elucidation of Capsular Polysaccha-
748 ride from *Streptococcus pneumoniae* Serogroup 16. *Journal of Bacteriology* **201**, 13 (2019).
- 749 [38] Lin, F. L. *et al.* Identification of the common antigenic determinant shared by *Streptococcus pneumoniae*
750 serotypes 33A, 35A, and 20 capsular polysaccharides. *Carbohydrate Research* **380**, 101–107 (2013). URL
751 <https://linkinghub.elsevier.com/retrieve/pii/S000862151300284X>.
- 752 [39] Lin, F. L. *et al.* Structure elucidation of capsular polysaccharides from *Streptococcus pneumoniae* serotype
753 33C, 33D, and revised structure of serotype 33B. *Carbohydrate Research* **383**, 97–104 (2014). URL
754 <https://linkinghub.elsevier.com/retrieve/pii/S0008621513003947>.
- 755 [40] Bush, C. A., Cisar, J. O. & Yang, J. Structures of Capsular Polysaccharide Serotypes 35F and 35C of
756 *Streptococcus pneumoniae* Determined by Nuclear Magnetic Resonance and Their Relation to Other Cross-
757 Reactive Serotypes. *Journal of Bacteriology* **197**, 2762–2769 (2015). URL <https://journals.asm.org/doi/10.1128/JB.00207-15>.

- 759 [41] Kondakova, A. N. *et al.* Reinvestigation of the O-antgens of Yersinia pseudotuberculosis: revision of the
760 O2c and confirmation of the O3 antigen structures. *Carbohydrate Research* **343**, 2486–2488 (2008). URL
761 <https://linkinghub.elsevier.com/retrieve/pii/S0008621508003443>.
- 762 [42] Henrissat, B. & Bairoch, A. Updating the sequence-based classification of glycosyl hydrolases. *The Bio-*
763 *chemical Journal* **316** (Pt 2), 695–696 (1996).
- 764 [43] Paton, J. C. & Trappetti, C. *Streptococcus pneumoniae* Capsular Polysaccharide. *Microbiology Spectrum*
765 **7**, 7.2.33 (2019). URL <https://journals.asm.org/doi/10.1128/microbiolspec.GPP3-0019-2018>.
- 766 [44] Doyle, L. *et al.* Mechanism and linkage specificities of the dual retaining β -Kdo glycosyltransferase modules
767 of KpsC from bacterial capsule biosynthesis. *Journal of Biological Chemistry* **299**, 104609 (2023). URL
768 <https://linkinghub.elsevier.com/retrieve/pii/S002192582300251X>.
- 769 [45] Krinos, C. M. *et al.* Extensive surface diversity of a commensal microorganism by multiple DNA inversions.
770 *Nature* **414**, 555–558 (2001). URL <https://www.nature.com/articles/35107092>.
- 771 [46] Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009). URL
772 <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-10-421>.
- 773 [47] Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide
774 sequences. *Bioinformatics* **22**, 1658–1659 (2006). URL <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btl158>.
- 775 [48] Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements
776 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>.
- 776 [49] Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching.
777 *Nucleic Acids Research* **39**, W29–W37 (2011). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr367>.
- 778 [50] Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a
779 multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009). URL
780 <https://academic.oup.com/bioinformatics/article/25/9/1189/203460>.
- 781 [51] Smith, T. & Waterman, M. Identification of common molecular subsequences. *Journal of Molecular Biology*
782 **147**, 195–197 (1981). URL <https://linkinghub.elsevier.com/retrieve/pii/0022283681900875>.
- 783 [52] 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>.
- 784 [53] Crippen, G. & Havel, T. *Distance Geometry and Molecular Conformation*, vol. 15 (Chemometrics Research
785 Studies Press Series, Research Studies Press, 1988).
- 786 [54] Huszczyński, S. M., Hao, Y., Lam, J. S. & Khursigara, C. M. Identification of the Pseudomonas aeruginosa
787 O17 and O15 O-Specific Antigen Biosynthesis Loci Reveals an ABC Transporter-Dependent Synthesis
788 Pathway and Mechanisms of Genetic Diversity. *Journal of Bacteriology* **202** (2020). URL <https://journals.asm.org/doi/10.1128/JB.00347-20>.
- 789 [55] Shannon, P. *et al.* Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction
790 Networks. *Genome Research* **13**, 2498–2504 (2003). URL <http://genome.cshlp.org/lookup/doi/10.1101/gr.1239303>.
- 791 [56] Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function using networkx
792 (2008). URL <https://www.osti.gov/biblio/960616>.
- 793 [57] Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep protein annotation. *BMC*
794 *Bioinformatics* **20**, 473 (2019). URL <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3019-7>.
- 795 [58] Petersen, B. O., Meier, S., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Determination of native
796 capsular polysaccharide structures of *Streptococcus pneumoniae* serotypes 39, 42, and 47F and comparison
797 to genetically or serologically related strains. *Carbohydrate Research* **395**, 38–46 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621514002560>.

- 807 [59] Petersen, B. O., Hindsgaul, O., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Structural elucidation
808 of the capsular polysaccharide from *Streptococcus pneumoniae* serotype 47A by NMR spectroscopy.
809 *Carbohydrate Research* **386**, 62–67 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621513004084>.
- 811 [60] Lam, J. S., Taylor, V. L., Islam, S. T., Hao, Y. & Kocíncová, D. Genetic and Functional Diversity of
812 *Pseudomonas aeruginosa* Lipopolysaccharide. *Frontiers in Microbiology* **2** (2011). URL <http://journal.frontiersin.org/article/10.3389/fmicb.2011.00118/abstract>.
- 814 [61] Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display
815 and annotation. *Nucleic Acids Research* **49**, W293–W296 (2021). URL <https://academic.oup.com/nar/article/49/W1/W293/6246398>.
- 817 [62] Hunter, J. D. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* **9**, 90–95
818 (2007). Publisher: IEEE COMPUTER SOC.
- 819 [63] Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nature Methods* **19**, 679–682 (2022).
820 URL <https://www.nature.com/articles/s41592-022-01488-1>.
- 821 [64] Schrödinger, LLC. The PyMOL Molecular Graphics System, Version 1.8 (2015).
- 822 [65] Shindyalov, I. N. & Bourne, P. E. Protein structure alignment by incremental combinatorial extension
823 (CE) of the optimal path. *Protein Engineering* **11**, 739–747 (1998).