

1 Introduction

2 (This paragraph can maybe be shortened a bit) Carbohydrate polymers (glycans) and glyco-conjugates are the
3 most abundant biomolecules on Earth and adopt a wide range of functions including energy storage, structure,
4 signaling, and mediators of host-pathogen interactions [1]. Due to the stereochemical diversity of monosaccharides
5 and the many possible linkages they can engage into, glycans display an enormous structural diversity
6 [2, 3]. Yet, our knowledge on their assembly is far from complete, especially in comparison to the enzymes
7 catalyzing their enzymatic breakdown.

8 The transfer of sugar moieties to acceptor molecules such as proteins, lipids or other sugars, is performed
9 by enzymes called glycosyltransferases or GTs [4]. Campbell and colleagues (1997) proposed a sequence-based
10 classification of GTs into 26 families, which was subsequently expanded to 65 families in 2003 [5]. The number of
11 sequence-based families has since continued to grow based on the necessary presence of at least one experimen-
12 tally characterized founding member to define a family. The constantly updated GT classification is presented
13 in the carbohydrate-active enzymes database (CAZy; www.cazy.org) along with similar family classifications
14 of other carbohydrate-active enzymes [6]. An additional advantage of the sequence-based classification is that
15 it readily enables genome mining for the presence of family members. Today there are 116 GT families in the
16 CAZy database and this number will continue growing as novel glycosyltransferases are progressively discovered
17 or as known GTs are incorporated in the database. In contrast to the EC numbers [7], the sequence-based
18 classification implicitly incorporates the structural features of GTs including the conservation of the catalytic
19 residues.

20 (This paragraph can maybe be shortened a bit) It was recognized very early that sequence-based GT families
21 group together enzymes that can utilize different sugar donors and/or acceptors, illustrating how GTs can evolve
22 to adopt novel substrates and form novel products [8, 5]. Mechanistically, glycosyltransferases can be either
23 retaining or inverting, based on the relative stereochemistry of the anomeric carbon of the sugar donor and of the
24 formed glycosidic bond [4]. This feature is conserved in previously defined sequence-based families, providing
25 predictive power to this classification, as the orientation of the glycosidic bond can be predicted safely even if
26 the precise transferred carbohydrate is not known.

27 The large majority of the 116 families of GTs listed in the CAZy database use donors activated by nucleotide
28 diphosphates. Eleven families utilize nucleotide monophospho-sugars (sialyl and KDO transferases), while 12
29 families utilize lipid monophospho-sugars. Only one family in the CAZy database utilizes lipid diphospho-
30 oligosaccharide donors: the oligosaccharyltransferases of family GT66, which transfer a pre-assembled oligosac-
31 charide to asparagine residues in N-glycoproteins [4, 9]. Several lipid diphospho-oligosaccharide-utilizing GTs
32 are currently missing in the CAZy database.

33 By contrast to the nucleotide-sugar dependent GTs, which are globular proteins with either a GT-A or
34 GT-B fold, the sugar-phospholipid-utilizing GTs are integral membrane proteins which have an overall GT-C
35 fold with a number of transmembrane helices that varies from 8 to 14 [4, 10]. Alexander and Locher recently
36 suggested two subgroups of GT-C glycosyltransferases, GT-C_A and GT-C_B, based on the structural features
37 of several of these families [10]. Several GTs from the GT-C_B subclass are missing in CAZy. They fall into
38 four major functional classes, which are all involved in the synthesis of bacterial cell wall polysaccharides and,
39 like CAZy family GT66, they catalyze the transfer of a glycan activated by the diphospholipid undecaprenyl
40 diphosphate (Und-PP).

41 The first functional group is the peptidoglycan polymerases, SEDS (shape, elongation, division and sporula-
42 tion) proteins. These proteins polymerize peptidoglycan in paris with class B penicillin-binding proteins, which
43 perform peptidoglycan crosslinking [11]. The structure of a SEDS protein from *Thermus thermophilus* has been
44 determined and consists of 10 transmembrane helices with several large extracellular loops containing function-
45 ally important residues [12]. A large hydrophobic groove containing highly conserved residues is thought to be
46 the lipid binding site. An Asp residue has been shown to be essential for RodA function in both *T. thermophilus*
47 and *B. subtilis* [12, 13].

48 The other three functional groups are involved in the synthesis of bacterial surface polysaccharides. Bacteria
49 synthesize various surface polysaccharides which confer them antigenic properties. Lipopolysaccharide (LPS)
50 is a polysaccharide specific of Gram-negative bacteria, and consists of the serotype-specific O-antigen attached
51 to the Lipid A-core oligosaccharide which is located in the outer membrane [14]. On the other hand capsular
52 polysaccharides (CPS, containing the K-antigen) are produced by both Gram-negative and Gram-positive
53 bacteria [15]. The covalent anchoring of CPS is still poorly understood, although it is found to be linked to
54 peptidoglycan in some Gram-positives [15]. Bacteria from the Enterobacteriales order produce yet another type
55 of surface polysaccharides referred to as the enterobacterial common antigen (ECA), which consists of repeating
56 units of N-acetylglucosamine, N-acetyl-D-mannosaminuronic acid and 4-acetamido-4,6-dideoxy-D-galactose [16].
57 Most of these surface polysaccharides are produced via the so-called Wzx/Wzy-dependent pathway, which takes
58 place on the plasma membrane (inner membrane in Gram-negatives) [17]. In this pathway, sugar repeat units
59 are assembled on an undecaprenyl-diphosphate (Und-PP) anchor on the cytoplasmic side of the membrane and
60 then flipped to the outside of the membrane by the flippase Wzx. The repeat units are then polymerized by the

bacterial polysaccharide polymerases (Wzy; BP-Pols), by transferring the growing polymer to the incoming new repeat units [17, 18]. In the case of LPS, the polymer (O-antigen) is then ligated onto Lipid A-core oligosaccharide by the O-antigen ligase (WaaL; O-Lig) [19]. ECA is produced via the same pathway, but with another set of enzymes including the polymerase (WzyE). In order to distinguish these polymerases from the serotype-specific polymerases, they are here referred to as enterobacterial common antigen polymerases (ECA-Pols).

In an attempt to complete the sequence-based classification of GTs, we have performed a detailed analysis of the primary sequence of SEDS proteins, ECA-Pols, BP-Pols and O-Ligs to assign their sequences to CAZy families and examined how sequence diversity correlates with the diversity of the transferred oligosaccharides and with the stereochemical outcome of the glycosyl transfer reaction.

2 Results

2.1 Peptidoglycan Polymerases

For building the CAZy family of SEDS proteins, we used the sequence from the published structure, 6BAR [12], as a starting point. The family GTxx1 was created and populated by using BLAST against Genbank, and subsequently using an HMM search against Genbank. GTxx1 is a very large family currently counting over 57,200 Genbank members in the CAZy database with a sequence similarity greater than 19% sequence identity over 221 residues.

The taxonomic distribution of family GTxx1 follows what was reported in [13], 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 [20] was used as seed sequence for the ECA-Pol family. Although the CAZy database only lists Genbank entries [21], we decided to build our multiple sequence alignments (MSAs) with 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. The ECA-Pols were assigned to a single new CAZy family, GTxx2. To date this new family contains over 4800 members with sequence identity greater than 38% over 414 residues, consistent with the conservation of acceptor, donor and product of the reaction.

As expected from their taxonomy-based designation, the ECA-Pol family (GTxx2) essentially contains sequences from the Enterobacterales order but also a few members of the Pasteurellales, suggesting that ECA-Pols of the latter were acquired by horizontal gene transfer (see below). why see below here?

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

2.3 O-antigen ligases

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

The greater diversity of the GTxx3 O-Ligs compared to the GTxx1 peptidoglycan polymerases and GTxx2 ECA-Pol appears in the form of the four divergent clades in the O-Lig phylogenetic tree (Supplementary Fig. 1). We hypothesize that this increased diversity originates from the extensive donor and moderate acceptor variability of O-Ligs [14]. Taxonomically, the GTxx3 O-Lig family is present in most bacteria, including both Gram-negatives and Gram-positives. The reaction performed by O-Ligs involves an inversion of the stereochemistry of the anomeric carbon since the sugar donor is axially bound to Und-PP and the reaction product is equatorially bound to Lipid A [19].

A recently discovered O-Lig, WadA, is bimodular with a GTxx3 domain appended to a globular glycosyltransferase domain of family GT25, which adds the last sugar to the oligosaccharide core [22]. We have constructed a tree with representative WadA homologs from the GTxx3 family (Supplementary Fig. 2) and observe that most of the sequences appended to a GT25 domain form one clade in the tree, except for a few outliers. This suggests a coupled action of the GT25 and of the GTxx3 at least for the bimodular O-ligs and

possibly for the entire family. The bimodular WadA O-Lig is observed in five genera including Mesorhizobium and Brucella.

2.4 Other bacterial polysaccharide polymerases

The fourth functional subgroup of GT-C_B are the BP-Pols. There is to our knowledge only one experimentally characterized BP-Pol [18]. However, several studies have identified BP-Pols from the polysaccharide gene clusters, and we decided to build our families based on these. We thus collected 363 predicted BP-Pol sequences (also referred to as Wzy) from seven review papers for various species, both Gram-negatives and Gram-positives: *Escherichia coli* [23], *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri* [24], *Salmonella enterica* [25], *Yersinia pseudotuberculosis*, *Yersinia similis* [26], *Pseudomonas aeruginosa* [17], *Acinetobacter baumanii*, *Acinetobacter nosocomialis* [27] and *Streptococcus pneumoniae* [28] (Supplementary Table 2).

In contrast to ECA-Pols, the donors as well as the acceptors of BP-Pols are highly variable. Others have reported an exceptional sequence diversity of BP-Pols even within the same species [17]. We also found that the sequences of BP-Pols are extremely diverse, and global alignments failed to reveal any conserved residue due to both sequence diversity and to the difficulty in aligning proteins with multiple and variable numbers of hydrophobic helices. It was therefore not possible to build a single family that could capture the diversity of BP-Pols.

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

Instead, we used a combination of SSN and HMM comparisons: First, we used an SSN with a strict threshold which would allow us to build good MSAs for the resulting clusters. This resulted in 204 clusters (Fig. 1a). Next, we created an HMM profile of each SSN cluster and compared the HMMs by all-vs-all pairwise HHblits, a program that aligns two HMMs and calculates a similarity score [29]. We then combined the SSN clusters into “superclusters” in a network analysis based on the HHblits scores (Fig. 1b). For this, we used a score cut-off of 160 in order to get a meaningful sequence and organismal diversity, resulting in 28 superclusters of varying sizes and 86 singleton clusters. Interestingly, the BP-Pols cluster across taxonomy, and even BP-Pols from Gram-positive and Gram-negative bacteria cluster together. The 14 largest superclusters have been included as new GT families in the CAZy database (GTxx4-GTx17) with a number of members ranging from 159 to 5,979 at 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 range of taxonomy, 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.

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 ([30, 24, 25, 26, 17, 27, 28]). Additionally, eight sugar structures were included, which were published after the review papers [31, 32, 33, 34, 35]. 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. **explain a bit more** (details in section 4.)

Having retrieved the sugar structures transferred by the BP-Pols, we first analyzed the stereochemistry of the bond catalyzed by the polymerase. As mentioned above, the stereochemical mechanism (inverting or retaining) is usually 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 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 seems likely that there is either an error in the chemical structure or that the serotype designation was incorrect.

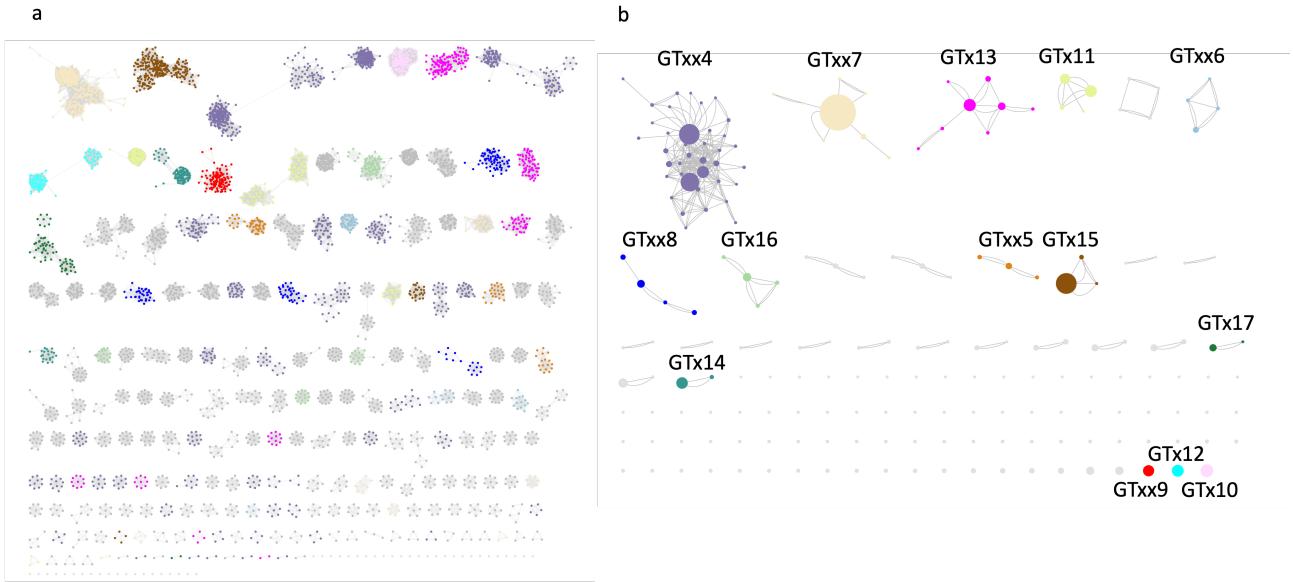


Figure 1: Clustering of BP-Pol sequences. a) SSN network with nodes representing proteins and edges representing pairwise alignment bit scores. b) HHblits network with nodes presenting SSN clusters and edges representing HHblits scores. The resulting clusters are referred to as “superclusters”. There 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 - GTxx17. In both a and b, the SSN clusters are coloured according to which supercluster they belong to.

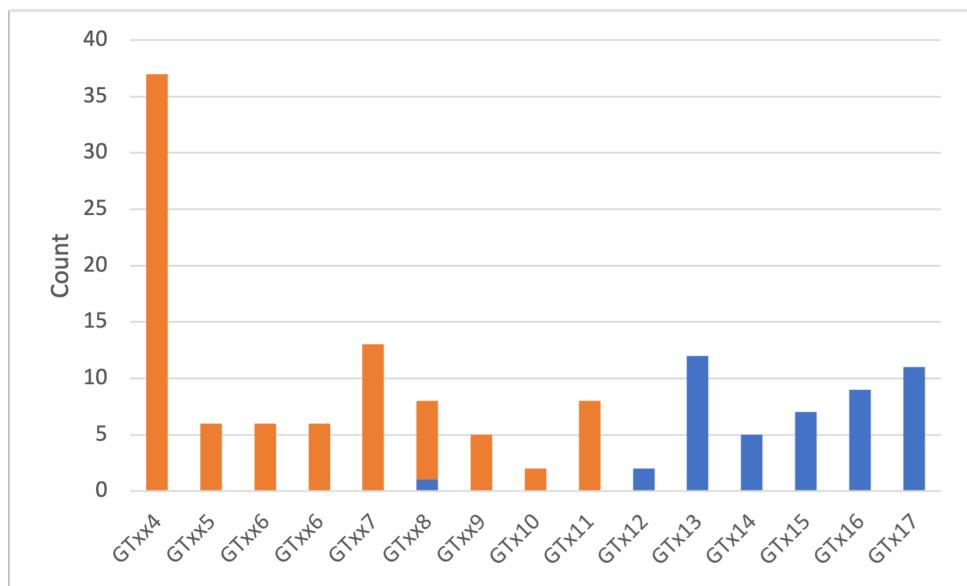


Figure 2: Conservation of stereochemical outcome of the reaction catalyzed in the various BP-Pol families. Equatorial bonds are shown in orange and axial bonds are shown in blue.

170 Next, we investigated whether there was a correlation between the structures of the transferred sugars
171 and the sequence similarity of the BP-Pols. We created phylogenetic trees of the BP-Pols in each family and
172 visualized them with the corresponding transferred repeat-units. We see that the sugars within each family
173 are similar, and correlates with the structure of the tree **reformulate** (Fig. 3, Supplementary Fig. 3). Notably,
174 we observe examples of BP-Pols from distant taxonomic serotypes that cluster in the same CAZy family and
175 have highly similar sugars. For example, (*Escherichia coli* O178 and *Streptococcus pneumoniae*) 47A in GTxx7
176 transfer sugars with almost identical backbones. There is only a slight variance in the middle of the repeat unit.
177 This suggest that horizontal gene transfer has occured. **reformulate**

178 To quantify the correlation between BP-Pol sequence and sugar structure, we developed an original pairwise
179 oligosaccharide similarity score. In our scoring scheme, the similarity of two glycans is estimated by examining
180 subsite moieties immediately upstream and downstream of the newly created interosidic bond, as we hypothesize
181 that these are the moieties most fitting the active site of the polymerase (Fig. 4). The minimum match between
182 two oligosaccharides corresponds to identical moieties at both subsites -1 and +1, which yields a score of 2.
183 Thereafter, the score increases by one unit for each additional match at contiguous subsites, -2, -3, etc., and
184 +2, +3, etc., up to a maximum value of 7 subsites found for the glycans encountered in this study (for details
185 see Methods).

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

189 2.6 Comparison of families

190 Others have previously reported sequence and structural similarity between RodA, O-Lig and some BP-Pols
191 [36, 10, 37, 13]. In order to investigate the relatedness of the new CAZy families, we compared the family
192 HMMs by all-vs-all HHblits analyses [29] (Fig. 6). Strikingly, we observe that the retaining BP-Pol families
193 cluster together on the heatmap together with the retaining ECA-Pols, while the inverting BP-Pols form two
194 distinct groups, one of them containing the inverting O-Ligs. The background noise between some inverting
195 and the retaining enzymes likely due to the general conservation of the successive transmembrane helices, which
196 is altered in the GTxx4-GTxx5-GTxx6 subgroup due to their different architecture (see below); on the other
197 hand, peptidoglycan polymerases segregate away from the other families.

198 In the CAZy database, clans have been defined for the glycoside hydrolases (GHs), which group together
199 CAZy families with distant sequence similarity, similar fold, similar catalytic machinery and stereochemical
200 outcome [38]. In extension of the report of the GT-C_B class by Alexander and Locher [10], and based on the
201 above-mentioned similarities between the new CAZy families, we can now define three clans within GT-C_B:
202 GT-C_{B1} consisting of inverting BP-Pol families and O-Lig, GT-C_{B2} consisting of retaining BP-Pol families and
203 ECA-Pol, and GT-C_{B3} consisting of inverting BP-Pol families (Table 1). The families within each clan share
204 residual, local, sequence similarity, insufficient to produce a multiple sequence alignment, but suggestive of
205 common ancestry.

206 In the absence of a three-dimensional structure, and based solely on the number of transmembrane helices,
207 we assigned clan GT-C_{B3} to the structural subclass GT-C_B of Alexander and Locher [10]. In addition, we also
208 present in Table 1 the families of GT-C glycosyltransferases that have not yet been assigned to a structural
209 class.

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

218 Most of the families in the inverting Clan GT-C_{B1} have two conserved Arg residues and one conserved either
219 Glu/Asp (in the BP-Pols) or His residue (in the O-Lig) (Fig. 7). In the pairwise HHblits alignments and
220 structural superimpositions, the Glu/Asp/His residues align, suggesting that they could play the same role. As
221 an example, the structural superimposition of the published O-Lig structure (7TPG) [37] and an AlphaFold
222 model from one representative of the inverting BP-Pol family GTxx8 is shown in Fig. 9a. The superimposition
223 produced an overall RMSD of 5.3 Å over 192 residues. Even with such a high RMSD, the two conserved Args
224 are oriented very similarly, and the conserved His in O-Lig is placed in the same position as the conserved Glu
225 in the BP-Pol. In O-Lig, the conserved His has been proposed to activate the acceptor, while the two Args are
226 proposed to position the donor by binding to the phosphate groups [37]. We hypothesize that the Glu and Asp
227 residues in the BP-Pols play the same role as the His in O-Lig.

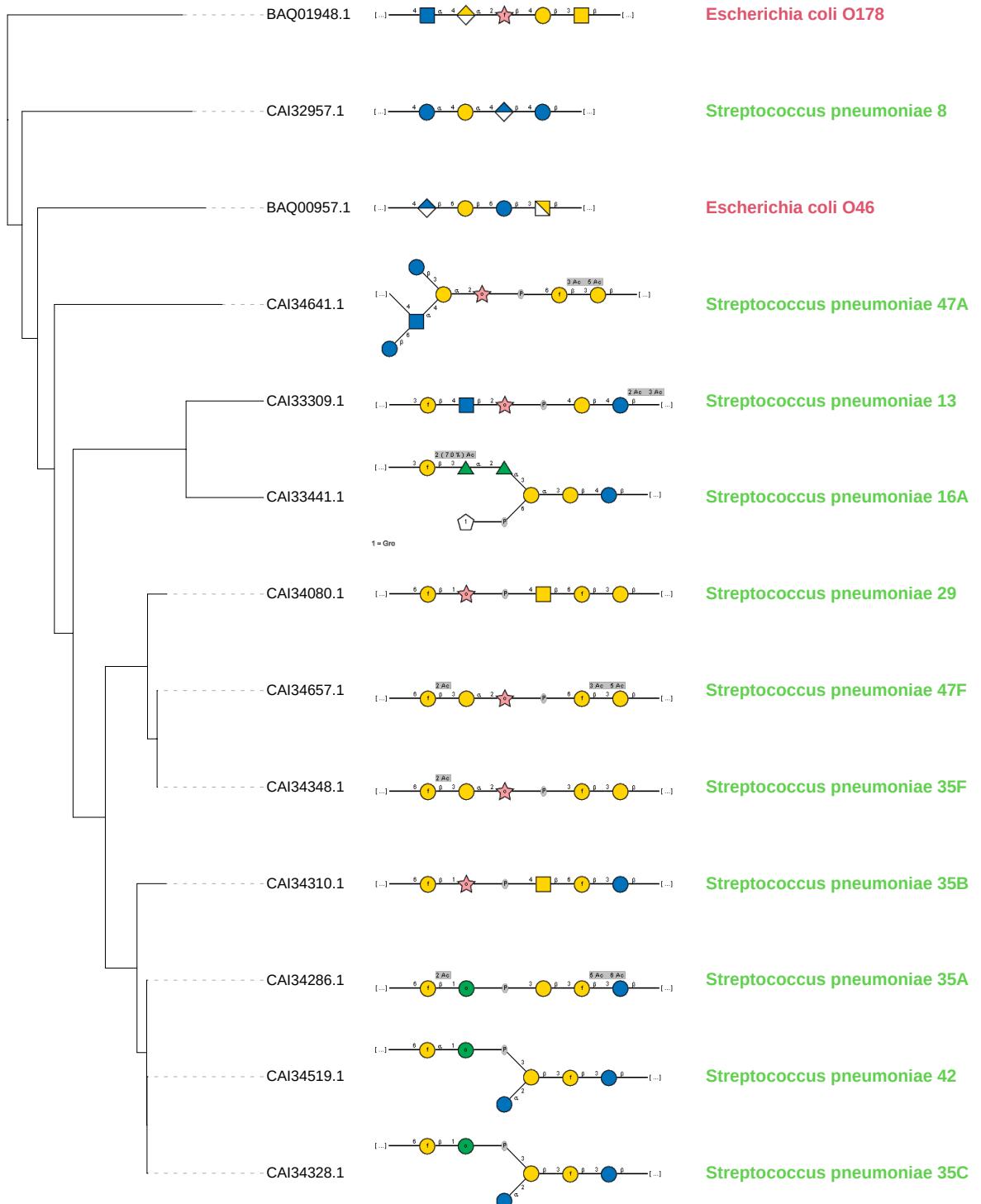


Figure 3: Phylogenetic tree of BP-Pols in family GTxx7 with structures of the corresponding sugar repeat units in SNFG format. The family contains BP-Pols from bacteria from distant taxonomies which transfer similar sugars. The trees for all the families are shown in Supplementary Figure 3.

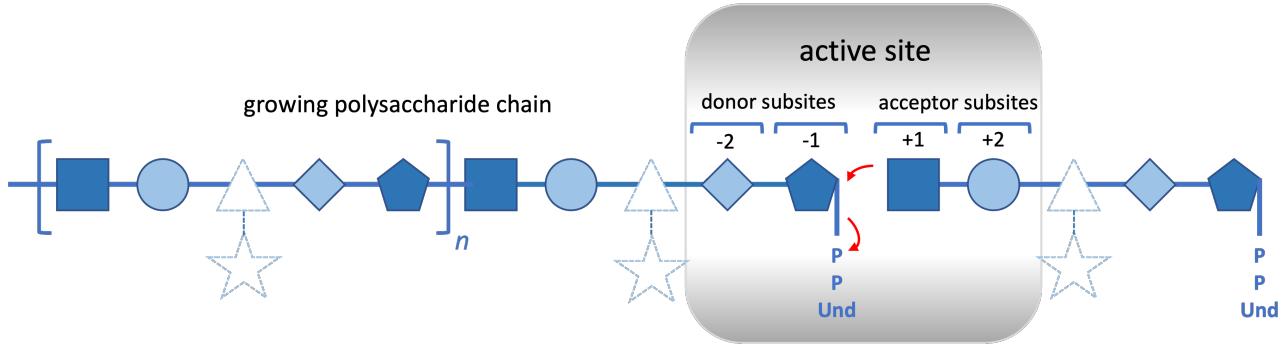


Figure 4: An idealized representation of a BP-Pol. The donor is the growing glycan chain activated by undecaprenyl pyrophosphate while the acceptor is a repeat unit monomer. The reaction is hypothesized to chiefly involve the sugar residues of the donor (subsites -2 and -1) and of the acceptor (subsites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

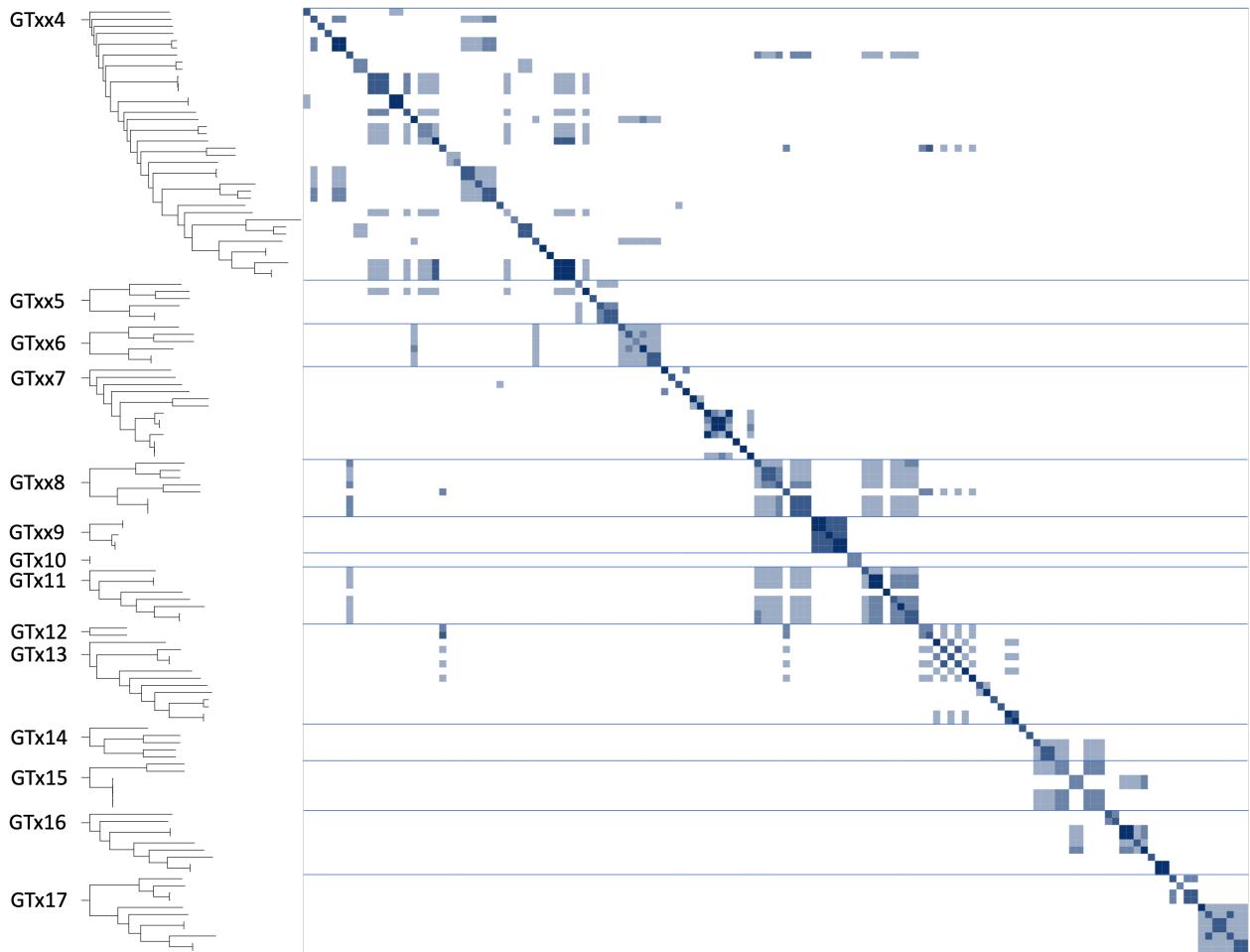


Figure 5: Glycan similarity of sugar repeat units polymerized by BP-Pols. All “seed” BP-Pols where the corresponding transferred oligosaccharide was known were included in the heatmap. A phylogenetic tree is shown for the polymerases in each CAZy family on the left. The glycan similarity scores are shown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

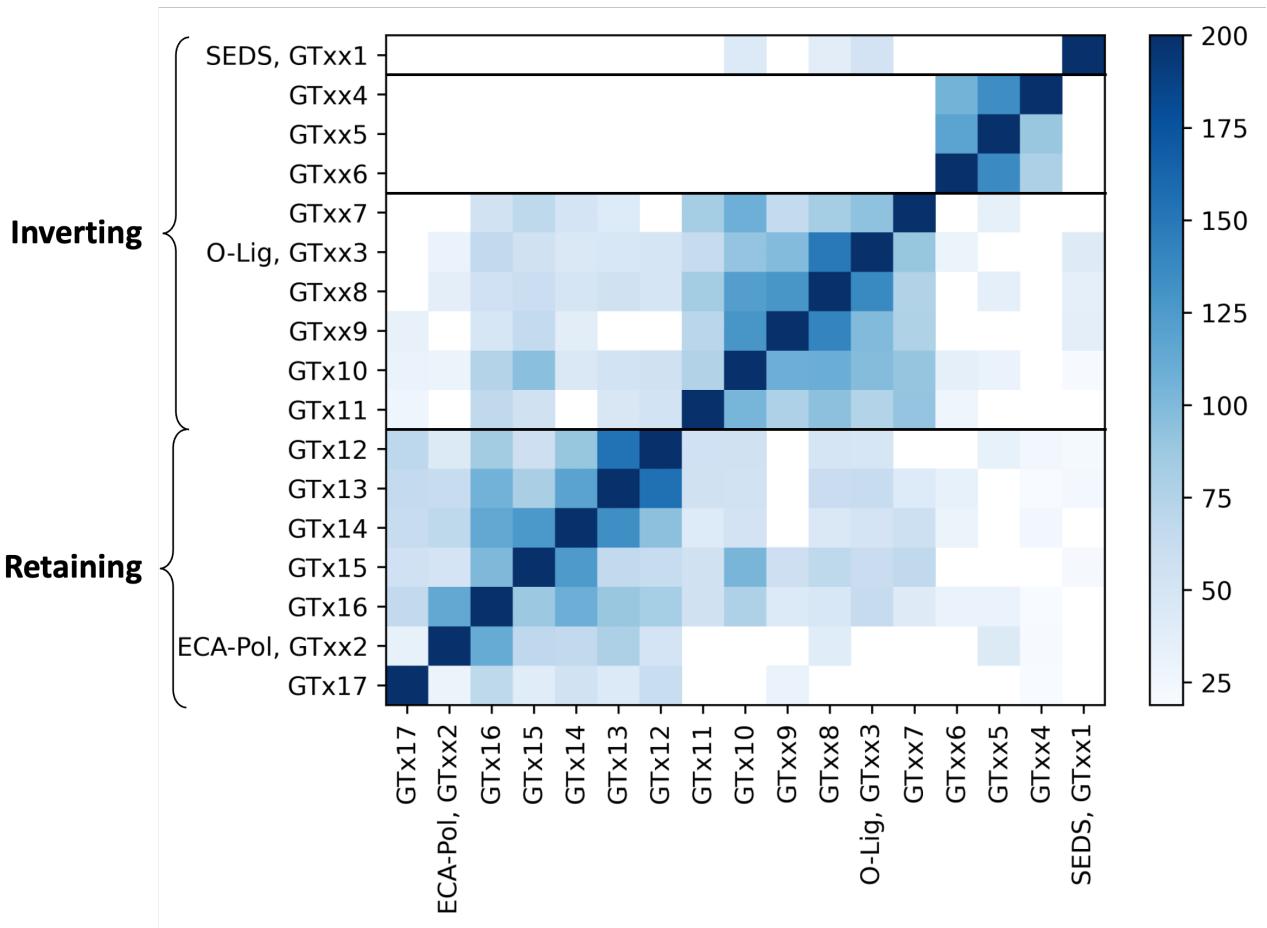


Figure 6: Heatmap of inter-family HHblits bit scores. The HHblits scores are shown on a color scale from white (low similarity score) to dark blue (high similarity score). The HHblits scores depend on the direction of the alignment, and therefore the heatmap is not symmetrical.

Structural subclass Alexander & Locher	CAZy clan	CAZy families	Mechanism	Donor
GT-CA (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

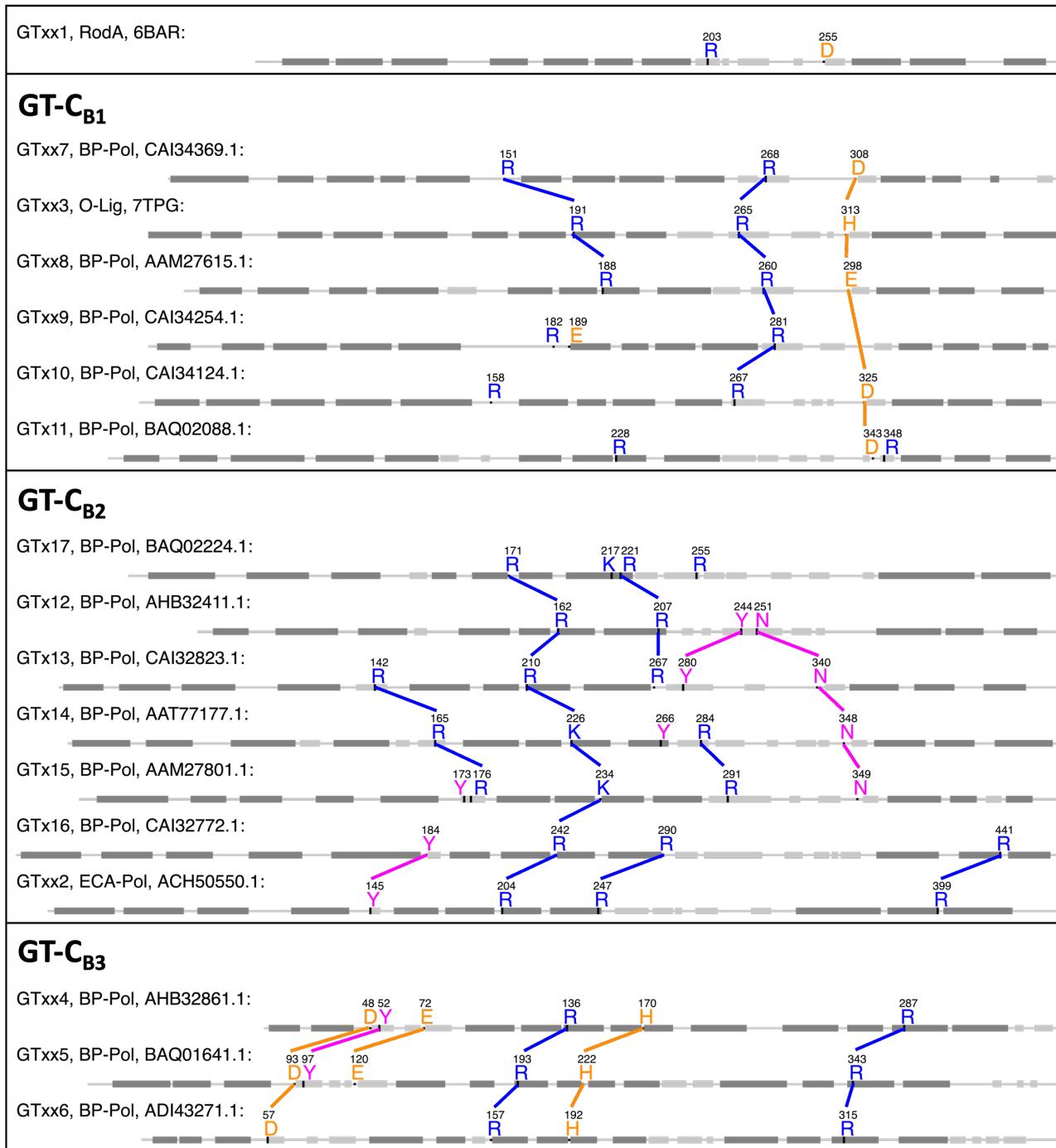
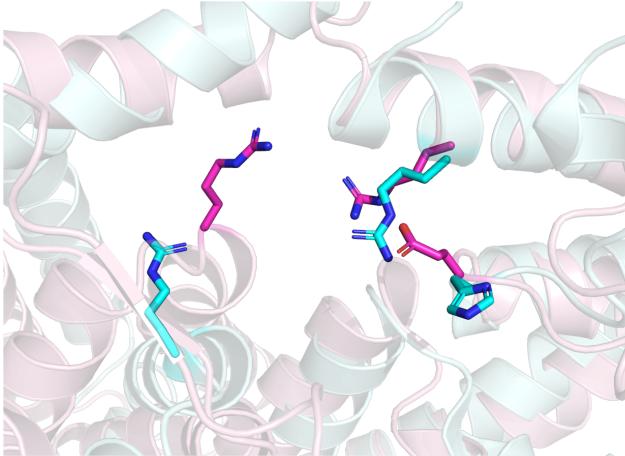
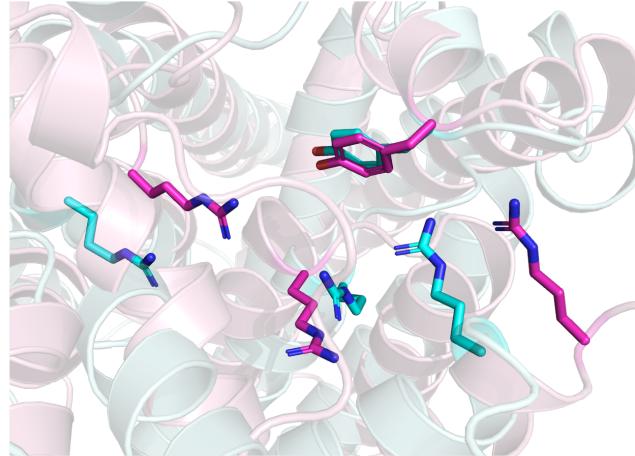


Figure 7: Comparison of conserved residues in the new GT families. The non-aliphatic conserved residues of each of the new CAZy families are shown on representative sequences. Transmembrane helices are shown in dark gray boxes, non-transmembrane helices are shown in light gray boxes. Lines are shown 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.



Cyan: O-Lig, GTxx3 (7TPG)
 Pink: BP-Pol, GTxx8 (AAM27615.1)



Cyan: ECA-Pol, GTxx2 (ACH50550.1)
 Pink: BP-Pol, GTx16 (CAI32772.1)

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 [37]. 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-C_{B2}, the pattern of conservation looks different. Here, most of the families have 2-3 conserved Arg/Lys and one conserved Tyr. As an example of the structural similarity in this clan, the structural superimposition of AlphaFold models from the ECA-Pol family GTxx2 and family GTx16 is shown in Fig 8b. The structures again show 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-C_{B2}.

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 [12]. 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 [37].

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 because of their extreme diversity even within a single species [23], 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, we were able to form larger families of similar polymerases from widely different taxonomies, thereby revealing conserved residues that are most likely functionally important.

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 [39]. Thus a hypothesis could be that the GTxx3 members in *S. pneumoniae* are responsible for the ligation of CPS to the peptidoglycan layer.

Because families are more robust when built with enough sequence diversity, many clusters of O-antigen polymerases were judged too small to build meaningful CAZy families. Additional polymerase families are thus

expected in the future with the accumulation of sequence data. For instance the small cluster that contains 47% identical BP-Pols from *E. coli* (GenBank BAQ01516.1) and *A. baumanii* (GenBank AHB32586.1) only contains eight sequences and will remain unclassified until enough sequence diversity has accumulated. This arbitrary decision comes from the need to devise a classification that can withstand a massive increase in the number of sequences without the need to constantly revise the content of the families. Thus new GT families based on O-antigen polymerases are poised to be formed when additional evidence becomes available.

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

It has been observed that for classical GT-A and GT-B fold glycosyltransferases, the catalytic mechanism is conserved within a family, but families with the same fold can have different mechanisms, possibly because the stereochemical outcome of the glycosyl transfer reaction is essentially dictated by the precise positioning and activation of the acceptor above (SN_2) or below (SN_i) the sugar ring of the donor [4]. Very occasionally, retaining glycosyltransferases have been shown to operate via a double displacement mechanism that involves Asp/Glu residues to form a glycosyl enzyme intermediate and to activate the acceptor that attacks this intermediate [40]. The families defined here display globally similar GT-C folds, and they also show conservation of the catalytic mechanism with about half of the families retaining and the other half inverting the anomeric configuration of the donor, suggesting that the outcome of the reaction catalyzed by GT-C glycosyltransferases is also dictated by the positioning of the acceptor with respect to the sugar plane of the acceptor. In turn this also suggests that retaining BP-Pols also operate by an SN_i mechanism rather than by the formation of a glycosyl enzyme intermediate. This hypothesis is supported by the lack of invariant Asp or Glu residues which could be involved in the formation and subsequent breakdown of a glycosyl enzyme intermediate in the retaining families GTxx2 and GTx12-GTx17. Additionally, the SN_i mechanism may provide protection against the interception of a glycosyl enzyme intermediate by a water molecule resulting in an undesirable hydrolysis reaction and termination of the polysaccharide elongation.

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

The analysis presented here shows that not only the stereochemistry of the glycosyl transfer is conserved in the BP-Pol families, but our development of an original method to estimate glycan similarity also reveals an unexpected degree of structural similarity of the oligosaccharide repeat units, suggesting that the latter constitutes a significant evolutionary constraint applying to the sequence and structure of BP-Pols. A closer inspection of the oligosaccharide repeat units within the families further reveals that the carbohydrates that appear the most constrained are the carbohydrates located (i) at the non-reducing end of the acceptor and (ii) close to the undecaprenyl pyrophosphate of the donor, i.e. the residues closest to the reaction center (Fig. 3). By contrast, residues away from the two extremities engaged in the polymerization reaction appear more variable, and can tolerate insertions/deletions or the presence of flexible residues such as linear glycerol or ribitol, with or without or the presence of a phosphodiester bond.

The version of the glycan similarity score presented here involves a direct translation of glycan IUPAC nomenclature into terms representing backbone configuration, i.e., ignoring chemical modifications and sidechains. Furthermore, a positive similarity score requires identical matches at both donor and acceptor positions (-1 and +1 sites in Fig. 3, respectively). These limitations will be addressed at a later stage (G.P. Gippert, in preparation).

We have next looked at the distribution of the new GT families in genomes, and particularly the families of BP-Pols. This uncovers broadly different schemes, with some bacteria having only one polymerase (and therefore only able to produce a single polysaccharide) while others having several, and sometimes more than 5, an observation in agreement with the report that *Bacteroides fragilis* produces no less than 8 different polysaccharides from distinct genomic loci [41]. The multiplicity of polysaccharide biosynthesis loci in some

322 genomes makes it sometimes difficult to assign a particular polysaccharide structure to a particular biosynthesis
323 operon. Although the families described here do not solve all problems, their correlation with the stereochemical
324 outcome of the glycosyl transfer reaction allows to resolve some inconsistencies (see above).

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

333 4 Methods

334 4.1 Alignment-based Clustering (Aclust)

335 Phylogenetic trees were generated using an in-house tool called Aclust (G.P.Gippert, manuscript in preparation)
336 comprising the following steps. (1) A distance matrix is computed from all-vs-all pairwise local pairwise
337 alignments [42], or from a multiple sequence alignment provided by MAFFT [43]. The distance calculation is
338 based on a variation of Scoredist ([44]), however with distance values normalized by sequence length rather
339 than alignment length. (2) The distance matrix is embedded into orthogonal coordinates using metric matrix
340 distance geometry [45], and a nearest-neighbor joining algorithm is used to create an initial tree. (3) Beginning
341 with the root node of the initial tree, each left and right subtree constitutes disjoint subsets of the original
342 sequence pool, which are embedded and rejoined separately (i.e., step 2 repeated for each subset), and the
343 process repeated recursively having the effect of gradually reducing deleterious effects on tree topology arising
344 from ‘long’ distances between unrelated proteins.

345 4.2 Building the peptidoglycan polymerase family (GTxx1)

346 The peptidoglycan polymerase family, GTxx1, was built by using “Blastp” from BLAST+ 2.12.0+ [46] against
347 Genbank with a threshold of approximately 30% to retrieve the family members. Next, an MSA was generated
348 with MAFFT v7.508 using the L-INS-i strategy (iterative refinement, using weighted sum-of-pairs and consis-
349 tency scores, of pairwise Needleman-Wunsch local alignments) [43] an HMM model was built using “hmmbuild”
350 from HMMER 3.3.2 [47]. The family was further populated using “hmmssearch” from HMMER 3.2.2 against
351 Genbank with a threshold of XX.

352 4.3 Building the Enterobacterial common antigen polymerases family (GTxx2)

353 A sequence library of ECA-Pols was constructed by using “Blastp” from BLAST+ 2.12.0+ [46] with the seed
354 sequence (Genbank accession AAC76800.1) against the NCBI non-redundant database version 61 with an E-
355 value threshold of 1e-60. The hits were redundancy reduced using CD-HIT 4.8.1 [48] with a threshold of 99%.
356 The redundancy reduced pool of ECA-Pol sequences was clustered using our in-house tool Aclust (see above),
357 and the tree showed one large clade and a few outliers. All the sequences in the large clade were used to build
358 an MSA using MAFFT v7.508 with the L-INS-i strategy [43]. An HMM was built based on this MSA using
359 the “hmmbuild” function from HMMER 3.3.2 [47]. The family GTxx3 was built in CAZy and populated using
360 “Blastp” against Genbank with an approximate threshold of 30% and hmmssearch against Genbank.

361 4.4 Building the O-antigen ligase family (GTxx3)

362 37 O-Lig sequences were selected from literature (Supplementary Table 1) and expanded using “blastp” against
363 the NCBI non-redundant database with an E-value cut-off of 1e-60. The blast hits were redundancy reduced
364 using CD-HIT with a threshold of 99%, resulting in a pool of 1,402 sequences. A phylogenetic tree of the
365 pool of O-Lig sequences was generated using Aclust (see section 4.1), which showed deep clefts between main
366 branches, and branches with sufficient internal diversity (Supplementary Figure 2). Based on these results, four
367 subfamilies were determined. An MSA was built for the family as well as for the subfamilies with MAFFT
368 v7.508 using the L-INS-i strategy. HMMs were built based on the MSAs using the “hmmbuild” function from
369 HMMER 3.3.2 [47]. The family was populated using Blastp against Genbank using an approximate threshold
370 of 30% identity with the seed sequences and using hmmssearch with the family and subfamily HMMs.

371 4.5 Building the Bacterial polysaccharide polymerase families (GTxx4-GTx17)

363 BP-Pol sequences were collected from review papers on biosynthesis of O-antigens and capsular polysaccharides in different species: *Escherichia coli* [23], *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri* [24], *Salmonella enterica* [25], *Yersinia pseudotuberculosis*, *Yersinia similis* [26], *Pseudomonas aeruginosa* [17], *Acinetobacter baumanii*, *Acinetobacter nosocomialis* [27] and *Streptococcus pneumoniae* [28] (complete list in Supplementary Table 2). The BP-Pols for *A. baumannii* O7 and O16 were omitted, because of uncertainty of their serotypes [27]. *P. aeruginosa* O15 was omitted, because it has been shown that the BP-Pol reported in [17] is inactivated and that the O-antigen is synthesized via the ABC-dependent pathway rather than the Wzx/Wzy-dependent pathway [49].

380 The sequence library was expanded using “blastp” 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 resulting in a pool of 20,850 sequences.

383 To find clusters of BP-Pol sequences that were big 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 [50] and MSAs of the members were built with MAFFT v7.508 using the L-INS-i strategy. The MSAs were inspected using Jalview [51], 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).

391 HMMs were then built for each SSN cluster using the “hmmbuild” function from HMMER 3.3.2, and the HMMs were compared using HHblits 3.3.0 [52]. 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 biggest 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 [53].

399 4.6 Analysis of sugar repeat-unit structures

400 In order to analyze the relation between BP-Pol sequence and structure of the transferred repeat-unit, we 401 retrieved the repeat-unit structures for the serotypes for the BP-Pols that were included in the new CAZy 402 families.

403 **Change this** The repeat-unit structures were retrieved in the following way:

- For *Acinetobacter baumannii*, the repeat unit structures of the O-antigens have been summarized in [27]. In this study, the linkage made by the polymerase was determined by analysis of the GTs in the gene clusters.
- For *Streptococcus pneumoniae*, the repeat unit structures of the capsular polysaccharides were summarized in [28], and the polymerase linkages were determined based on the initial transferase and the other GTs in the polysaccharide gene cluster. Eight additional repeat-unit structures were included which were elucidated after the review paper; from *S. pneumoniae* 16A [31], 33A [32], 33C and 33D [33], 35C and 35F [34], 42 and 47F [54] and 47A [55]. The polymerase linkage has been determined in all these studies, except for those of *S. pneumoniae* 33A and 47A. For 33A, we determined the polymerase linkage based on the gene cluster having the initial transferase WchA, which transfers a glucose [27]. 47A has WcjG as the initial transferase, which transfers Gal or Galp [27]. Since the repeat-unit contains both Gal and Galp, we could not determine the polymerase linkage unambiguously. However, the sugar unit is very similar to other sugars in the family (most similar to *S. pneumoniae* 13, and we proposed the equivalent phase of that one. Finally, a revised structure has been published of 33B [33], and we used that structure instead.
- For *Yersinia pseudotuberculosis*, the repeat unit structures of the O-antigens have been summarized in [26], in which the polymerase linkages were also determined. For *Y. pseudotuberculosis* O3, we used the revised structure [35].
- For *Salmonella enterica*, the repeat unit structures of the O-antigens have been summarized in [25]. In this study, the linkage made by the polymerase was determined by analysis of the GTs in the gene clusters.
- For *Escherichia coli*, the repeat unit structures of the O-antigens have been summarized in [30]. **check how they determined the polymerase linkage**

- 425 • For *Shigella*, the repeat unit structures of the O-antigens have been summarized in [24]. In this study,
 426 the linkage made by the polymerase was determined based on the initial sugar unit being GlcNAc or
 427 GalpNAc.
- 428 • For *Pseudomonas aeruginosa*, the repeat unit structures were retrieved from the review paper [17]. check
 429 how they determined phases *Pseudomonas aeruginosa* O2 and O16 contain two BP-Pol genes; one BP-Pol
 430 localized in the O-antigen biosynthesis cluster, which polymerizes the sugar repeat units with an α bond
 431 and one BP-Pol localized outside the biosynthesis cluster which polymerizes the repeat units with a β
 432 bond [56]. The polymerases in our dataset are the ones that polymerize the α bond and we therefore
 433 report the sugar structure with the alpha bond.

434 The CSDB database (<http://csdb.glycoscience.ru>) [57] was used for finding literature and retrieving linear
 435 sugar strings and SNFG image representations of the repeat-unit structures.

436 Phylogenetic trees for only seed sequences in each of the newly created BP-Pol families were generated
 437 using MAFFT v7.508 [43] to supply an initial multiple sequence alignment, followed by Aclust (section 4.2) for
 438 distance matrix embedding and clustering. Seed sequences are those where the sugar repeat-unit structure is
 439 known. The trees were visualized with the corresponding sugar structures in iTOL [58].

440 4.7 Oligosaccharide backbone similarity score

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

447 To facilitate the scoring, we have chosen to first translate oligosaccharides from IUPAC nomenclature into
 448 a set of simplified geometric subunits that represent only monomer dimension and stereochemistry of acceptor
 449 and anomeric donor carbon atoms, thus focusing entirely on the glycan backbone (Fig. 9). Briefly, the monomer
 450 dimension is represented by a single letter P, F or L depending on whether the monomer sugar is a pyranose,
 451 furanose or linear/open, respectively. Stereochemistry of the acceptor and donor carbon atoms is represented
 452 by the index number of the carbon position within the ring/monomer, followed by a single letter U, D or
 453 N depending on whether the linked oxygen atom is U (up=above) the monomer ring, D (down=below) the
 454 monomer ring, or N (neither above or below the ring), this latter category is assigned in the cases of extensive
 455 conformational flexibility such as with alditols or C6 linkages. Chemical modifications, side chains, and the
 456 configuration of non-linking carbons are ignored. Further details, limitations and extensions will be presented
 457 elsewhere (G.P. Gippert, manuscript in preparation).

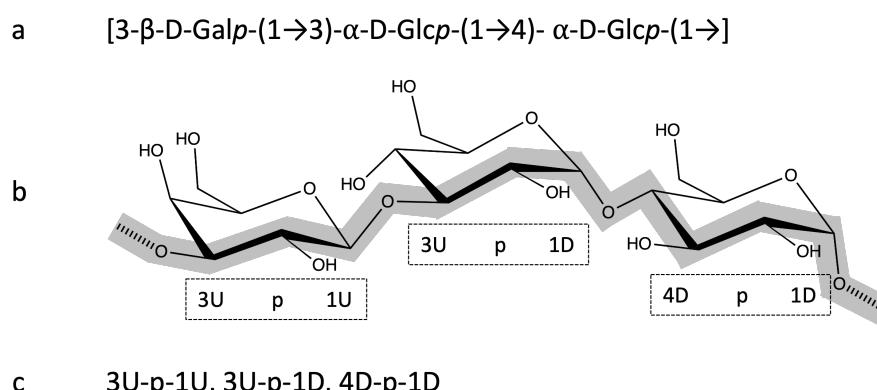


Figure 9: Oligosaccharide translation from IUPAC nomenclature to backbone (geometric) subunits for a trisaccharide consisting of one D-galactopyranose and two D-glucopyranose residues joined by intramolecular $\beta 1\rightarrow 3$ and $\alpha 1\rightarrow 4$ bonds, respectively, and an intermolecular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

458 4.8 Comparison of the families

459 Pairwise HHblits analyses [29] were performed for each of the new CAZy families. The HHblits scores were
460 visualized in a heatmap using Python Matplotlib [59].

461 AlphaFold2 [13] structures were generated of representative proteins from the families using the Colab-
462 Fold implementation [60] on our internal GPU cluster processed with the recommended settings. The best
463 ranked relaxed model was used. The protein structures were visualized in PyMOL [61] and pairwise structural
464 superimpositions were performed using the CEalign algorithm [62].

465 5 Data availability

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

469 6 Acknowledgements

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

473 7 Author contributions

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

478 8 Competing interests

479 None

480 References

- 481 [1] Varki, A. *et al.* (eds.) *Essentials of Glycobiology* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor
482 (NY), 2022), 4th edn. URL <http://www.ncbi.nlm.nih.gov/books/NBK579918/>.
- 483 [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
484 saccharide sequencing or synthesis systems. *Glycobiology* **4**, 759–767 (1994).
- 486 [3] Lapébie, P., Lombard, V., Drula, E., Terrapon, N. & Henrissat, B. Bacteroidetes use thousands of enzyme
487 combinations to break down glycans. *Nature Communications* **10**, 2043 (2019). URL <https://www.nature.com/articles/s41467-019-10068-5>.
- 489 [4] Lairson, L., Henrissat, B., Davies, G. & Withers, S. Glycosyltransferases: Structures, Functions, and
490 Mechanisms. *Annual Review of Biochemistry* **77**, 521–555 (2008). URL <https://www.annualreviews.org/doi/10.1146/annurev.biochem.76.061005.092322>.
- 492 [5] Coutinho, P. M., Deleury, E., Davies, G. J. & Henrissat, B. An Evolving Hierarchical Family Classification
493 for Glycosyltransferases. *Journal of Molecular Biology* **328**, 307–317 (2003). URL <https://linkinghub.elsevier.com/retrieve/pii/S0022283603003073>.
- 495 [6] Drula, E. *et al.* The carbohydrate-active enzyme database: functions and literature. *Nucleic Acids Research*
496 **50**, D571–D577 (2022). URL <https://academic.oup.com/nar/article/50/D1/D571/6445960>.
- 497 [7] McDonald, A. G. & Tipton, K. F. Fifty-five years of enzyme classification: advances and difficulties. *The FEBS journal* **281**, 583–592 (2014).

- 499 [8] Campbell, J. A., Davies, G. J., Bulone, V. & Henrissat, B. A classification of nucleotide-diphospho-sugar
500 glycosyltransferases based on amino acid sequence similarities. *The Biochemical Journal* **326** (Pt 3),
501 929–939 (1997).
- 502 [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>.
- 503 [10] 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>.
- 504 [11] 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>.
- 505 [12] 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>.
- 506 [13] Meeske, A. J. *et al.* SEDS proteins are a widespread family of bacterial cell wall polymerases. *Nature* **537**,
507 634–638 (2016). URL <http://www.nature.com/articles/nature19331>.
- 508 [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>.
- 509 [15] 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>.
- 510 [16] 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>.
- 511 [17] 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>.
- 512 [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>.
- 513 [19] 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>.
- 514 [20] 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>.
- 515 [21] 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>.
- 516 [22] 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>.
- 517 [23] 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>.
- 518 [24] 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>.
- 519 [25] 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>.

- 547 [26] Kenyon, J. J., Cunneen, M. M. & Reeves, P. R. Genetics and evolution of *Yersinia pseudotuberculosis* O-
548 specific polysaccharides: a novel pattern of O-antigen diversity. *FEMS Microbiology Reviews* **41**, 200–217
549 (2017). URL <https://academic.oup.com/femsre/article/41/2/200/2996588>.
- 550 [27] Hu, D., Liu, B., Dijkshoorn, L., Wang, L. & Reeves, P. R. Diversity in the Major Polysaccharide Antigen
551 of *Acinetobacter baumannii* Assessed by DNA Sequencing, and Development of a Molecular Serotyping
552 Scheme. *PLoS ONE* **8**, e70329 (2013). URL <https://dx.plos.org/10.1371/journal.pone.0070329>.
- 553 [28] Bentley, S. D. *et al.* Genetic Analysis of the Capsular Biosynthetic Locus from All 90 Pneumococcal
554 Serotypes. *PLoS Genetics* **2**, e31 (2006). URL <https://dx.plos.org/10.1371/journal.pgen.0020031>.
- 555 [29] Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence
556 searching by HMM-HMM alignment. *Nature Methods* **9**, 173–175 (2012). URL <http://www.nature.com/articles/nmeth.1818>.
- 558 [30] Liu, B. *et al.* Structure and genetics of *Escherichia coli* O antigens. *FEMS Microbiology Reviews* **44**,
559 655–683 (2020). URL <https://academic.oup.com/femsre/article/44/6/655/5645236>.
- 560 [31] Li, C. *et al.* Structural, Biosynthetic, and Serological Cross-Reactive Elucidation of Capsular Polysaccha-
561rides from *Streptococcus pneumoniae* Serogroup 16. *Journal of Bacteriology* **201**, 13 (2019).
- 562 [32] Lin, F. L. *et al.* Identification of the common antigenic determinant shared by *Streptococcus pneumoniae* serotypes
563 33A, 35A, and 20 capsular polysaccharides. *Carbohydrate Research* **380**, 101–107 (2013). URL
564 <https://linkinghub.elsevier.com/retrieve/pii/S000862151300284X>.
- 565 [33] Lin, F. L. *et al.* Structure elucidation of capsular polysaccharides from *Streptococcus pneumoniae* serotype
566 33C, 33D, and revised structure of serotype 33B. *Carbohydrate Research* **383**, 97–104 (2014). URL
567 <https://linkinghub.elsevier.com/retrieve/pii/S0008621513003947>.
- 568 [34] Bush, C. A., Cisar, J. O. & Yang, J. Structures of Capsular Polysaccharide Serotypes 35F and 35C of
569 *Streptococcus pneumoniae* Determined by Nuclear Magnetic Resonance and Their Relation to Other Cross-
570 Reactive Serotypes. *Journal of Bacteriology* **197**, 2762–2769 (2015). URL <https://journals.asm.org/doi/10.1128/JB.00207-15>.
- 572 [35] Kondakova, A. N. *et al.* Reinvestigation of the O-antigens of *Yersinia pseudotuberculosis*: revision of the
573 O2c and confirmation of the O3 antigen structures. *Carbohydrate Research* **343**, 2486–2488 (2008). URL
574 <https://linkinghub.elsevier.com/retrieve/pii/S0008621508003443>.
- 575 [36] Nygaard, R. *et al.* Structural basis of peptidoglycan synthesis by *E. coli* RodA-PBP2 complex. *Nature
576 Communications* **14**, 5151 (2023). URL <https://www.nature.com/articles/s41467-023-40483-8>.
- 577 [37] Ashraf, K. U. *et al.* Structural basis of lipopolysaccharide maturation by the O-antigen ligase. *Nature* **604**,
578 371–376 (2022). URL <https://www.nature.com/articles/s41586-022-04555-x>.
- 579 [38] Henrissat, B. & Bairoch, A. Updating the sequence-based classification of glycosyl hydrolases. *The Bio-
580 chemical Journal* **316** (Pt 2), 695–696 (1996).
- 581 [39] Paton, J. C. & Trappetti, C. *Streptococcus pneumoniae* Capsular Polysaccharide. *Microbiology Spectrum*
582 **7**, 7.2.33 (2019). URL <https://journals.asm.org/doi/10.1128/microbiolspec.GPP3-0019-2018>.
- 583 [40] Doyle, L. *et al.* Mechanism and linkage specificities of the dual retaining β-Kdo glycosyltransferase modules
584 of KpsC from bacterial capsule biosynthesis. *Journal of Biological Chemistry* **299**, 104609 (2023). URL
585 <https://linkinghub.elsevier.com/retrieve/pii/S002192582300251X>.
- 586 [41] Krinos, C. M. *et al.* Extensive surface diversity of a commensal microorganism by multiple DNA inversions.
587 *Nature* **414**, 555–558 (2001). URL <https://www.nature.com/articles/35107092>.
- 588 [42] Smith, T. & Waterman, M. Identification of common molecular subsequences. *Journal of Molecular Biology*
589 **147**, 195–197 (1981). URL <https://linkinghub.elsevier.com/retrieve/pii/0022283681900875>.
- 590 [43] Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements
591 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>.
- 593 [44] 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>.

- 595 [45] Crippen, G. & Havel, T. *Distance Geometry and Molecular Conformation*, vol. 15 (Chemometrics Research
596 Studies Press Series, Research Studies Press, 1988).
- 597 [46] Camacho, C. *et al.* BLAST+: architecture and applications. *BMC Bioinformatics* **10**, 421 (2009). URL
598 <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/1471-2105-10-421>.
- 599 [47] Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching.
600 *Nucleic Acids Research* **39**, W29–W37 (2011). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkr367>.
- 602 [48] Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide
603 sequences. *Bioinformatics* **22**, 1658–1659 (2006). URL <https://academic.oup.com/bioinformatics/article-lookup/doi/10.1093/bioinformatics/btl158>.
- 605 [49] Huszcynski, S. M., Hao, Y., Lam, J. S. & Khursigara, C. M. Identification of the Pseudomonas aeruginosa
606 O17 and O15 O-Specific Antigen Biosynthesis Loci Reveals an ABC Transporter-Dependent Synthesis
607 Pathway and Mechanisms of Genetic Diversity. *Journal of Bacteriology* **202** (2020). URL <https://journals.asm.org/doi/10.1128/JB.00347-20>.
- 609 [50] Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function using networkx
610 (2008). URL <https://www.osti.gov/biblio/960616>.
- 611 [51] Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a
612 multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009). URL
613 <https://academic.oup.com/bioinformatics/article/25/9/1189/203460>.
- 614 [52] Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep protein annotation. *BMC
615 Bioinformatics* **20**, 473 (2019). URL <https://bmcbioinformatics.biomedcentral.com/articles/10.1186/s12859-019-3019-7>.
- 617 [53] Shannon, P. *et al.* Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction
618 Networks. *Genome Research* **13**, 2498–2504 (2003). URL <http://genome.cshlp.org/lookup/doi/10.1101/gr.1239303>.
- 620 [54] Petersen, B. O., Meier, S., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Determination of native
621 capsular polysaccharide structures of Streptococcus pneumoniae serotypes 39, 42, and 47F and comparison
622 to genetically or serologically related strains. *Carbohydrate Research* **395**, 38–46 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621514002560>.
- 624 [55] Petersen, B. O., Hindsgaul, O., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Structural elucidation
625 of the capsular polysaccharide from Streptococcus pneumoniae serotype 47A by NMR spectroscopy.
626 *Carbohydrate Research* **386**, 62–67 (2014). URL <https://linkinghub.elsevier.com/retrieve/pii/S0008621513004084>.
- 628 [56] Lam, J. S., Taylor, V. L., Islam, S. T., Hao, Y. & Kocíncová, D. Genetic and Functional Diversity of
629 Pseudomonas aeruginosa Lipopolysaccharide. *Frontiers in Microbiology* **2** (2011). URL <http://journal.frontiersin.org/article/10.3389/fmicb.2011.00118/abstract>.
- 633 [57] Toukach, P. V. & Egorova, K. S. Carbohydrate structure database merged from bacterial, archaeal, plant
and fungal parts. *Nucleic Acids Research* **44**, D1229–D1236 (2016). URL <https://academic.oup.com/nar/article-lookup/doi/10.1093/nar/gkv840>.
- 637 [58] Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display
638 and annotation. *Nucleic Acids Research* **49**, W293–W296 (2021). URL <https://academic.oup.com/nar/article/49/W1/W293/6246398>.
- 642 [59] Hunter, J. D. Matplotlib: A 2D graphics environment. *Computing in Science & Engineering* **9**, 90–95
643 (2007). Publisher: IEEE COMPUTER SOC.
- 647 [60] Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nature Methods* **19**, 679–682 (2022).
URL <https://www.nature.com/articles/s41592-022-01488-1>.
- 651 [61] Schrödinger, LLC. The PyMOL Molecular Graphics System, Version 1.8 (2015).
- 655 [62] Shindyalov, I. N. & Bourne, P. E. Protein structure alignment by incremental combinatorial extension
656 (CE) of the optimal path. *Protein Engineering* **11**, 739–747 (1998).