

<https://doi.org/10.1038/s42003-024-05930-2>

# Diversity of sugar-diphospholipid-utilizing glycosyltransferase families

Ida K. S. Meitil <sup>1</sup>, Garry P. Gippert<sup>1</sup>, Kristian Barrett <sup>1</sup>, Cameron J. Hunt <sup>1</sup> & Bernard Henrissat <sup>1,2,3</sup> ✉

Peptidoglycan polymerases, enterobacterial common antigen polymerases, O-antigen ligases, and other bacterial polysaccharide polymerases (BP-Pols) are glycosyltransferases (GTs) that build bacterial surface polysaccharides. These integral membrane enzymes share the particularity of using diphospholipid-activated sugars and were previously missing in the carbohydrate-active enzymes database (CAZy; [www.cazy.org](http://www.cazy.org)). While the first three classes formed well-defined families of similar proteins, the sequences of BP-Pols were so diverse that a single family could not be built. To address this, we developed a new clustering method using a combination of a sequence similarity network and hidden Markov model comparisons. Overall, we have defined 17 new GT families including 14 of BP-Pols. We find that the reaction stereochemistry appears to be conserved in each of the defined BP-Pol families, and that the BP-Pols within the families transfer similar sugars even across Gram-negative and Gram-positive bacteria. Comparison of the new GT families reveals three clans of distantly related families, which also conserve the reaction stereochemistry.

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

The transfer of sugar moieties to acceptor molecules such as proteins, lipids or other sugars, is catalyzed by enzymes called glycosyltransferases or GTs<sup>4</sup>. Campbell and colleagues proposed a sequence-based classification of GTs into 26 families<sup>5</sup>. The number of sequence-based families has since continued to grow based on the necessary presence of at least one experimentally characterized founding member to define a family, and is presented in the carbohydrate-active enzymes database (CAZy; [www.cazy.org](http://www.cazy.org))<sup>6</sup>. An advantage of the sequence-based classification is that it readily enables genome mining for the presence of new family members. Today there are 118 GT families in the CAZy database and in contrast to the EC numbers<sup>7</sup>, the sequence-based classification implicitly incorporates the structural features of GTs including the conservation of the catalytic residues.

It was recognized very early that sequence-based GT families group together enzymes that can utilize different sugar donors and/or acceptors,

illustrating how GTs can evolve to adopt novel substrates and form novel products<sup>5,8</sup>. Mechanistically, glycosyltransferases can be either retaining or inverting, based on the relative stereochemistry of the anomeric carbon of the sugar donor and of the formed glycosidic bond<sup>4</sup>. With almost no exceptions, this feature is conserved in previously defined sequence-based families, providing predictive power to this classification, as the orientation of the glycosidic bond can be predicted even if the precise transferred carbohydrate is not known.

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

The first of these four functional classes corresponds to the peptidoglycan polymerases, shape, elongation, division and sporulation (SEDS) proteins. These proteins polymerize peptidoglycan in complex with class B penicillin-binding proteins<sup>10</sup>. Several 3-D structures of SEDS proteins have

<sup>1</sup>Department of Biotechnology and Biomedicine, Technical University of Denmark, Kgs. Lyngby, Denmark. <sup>2</sup>Architecture et Fonction des Macromolécules Biologiques, CNRS, Aix-Marseille Université, Marseille, France. <sup>3</sup>Department of Biological Sciences, King Abdulaziz University, Jeddah, Saudi Arabia.

✉ e-mail: [bernard.henrissat@gmail.com](mailto:bernard.henrissat@gmail.com)

been determined, and they harbor 10 transmembrane helices and one long extracellular loop<sup>11–13</sup>. This loop contains an Asp residue, which has been shown to be essential for SEDS function<sup>11,14</sup>.

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

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

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

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

## Results

### Peptidoglycan polymerases

For building the CAZy family of SEDS proteins, we used four characterized proteins as seed sequences: the proteins with PDB IDs 6BAR<sup>11</sup>, 8TJ3<sup>13</sup> and 8BH1<sup>12</sup>, and the protein with GenBank accession CAB15838.1<sup>24</sup>. Family GT119 was created and initially populated by using BLAST against GenBank, and subsequently by searching against GenBank with a hidden

Markov model (HMM) built from the retrieved sequences. GT119 is a very large family currently counting over 57,200 GenBank members in the CAZy database with a pairwise sequence identity of 19% over 221 residues for the most distant members. The taxonomic distribution of family GT119 follows what was reported in<sup>14</sup>, 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 Und-PP is  $\alpha$ -linked. The carbohydrate repeat unit of peptidoglycan being  $\beta$ -linked, the glycosyl transfer reaction thus inverts the stereochemistry of the anomeric carbon involved in the newly formed glycosidic bond.

### Enterobacterial common antigen polymerases

The ECA-Pol which was studied by Maczuga et al.<sup>25</sup> was used as seed sequence for building the ECA-Pol family. Although the CAZy database only lists GenBank entries<sup>26</sup>, we decided to build our multiple sequence alignments (MSAs) with sequences from the NCBI non-redundant database in order to capture more diversity. An ECA-Pol sequence library was thus constructed from the seed sequence using BLAST against the non-redundant database of the NCBI. The ECA-Pols were assigned to a single new CAZy family, GT120. To date this new family contains over 4800 GenBank members with high similarity (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 (GT120) 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. 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.

### O-antigen ligases

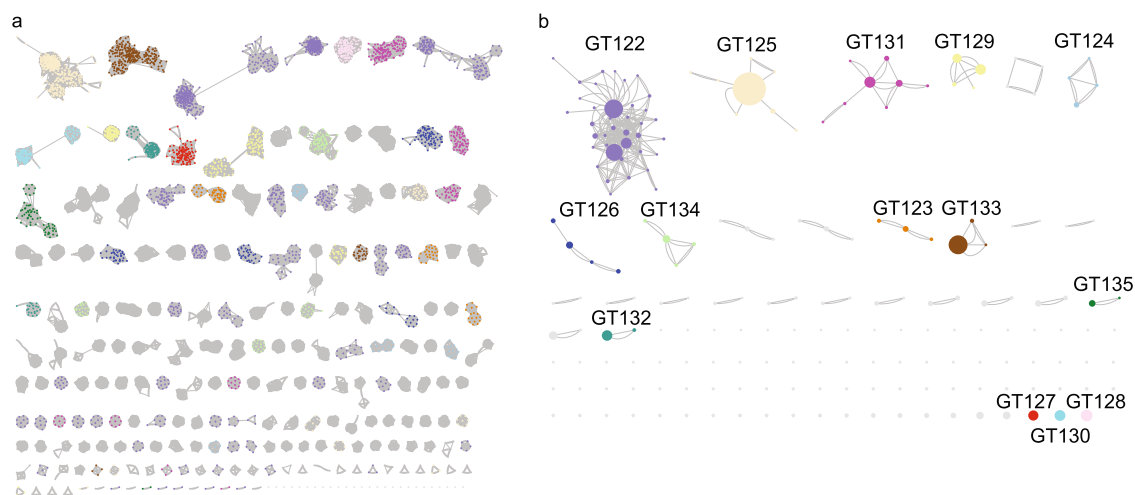
With the aim of including the O-Ligs in the CAZy database, we collected 37 O-Lig sequences (Supplementary Data 1) and constructed a sequence library from these seed sequences using BLAST against the NCBI non-redundant database. A phylogenetic tree of the sequence library revealed four distantly related clades (Supplementary Fig. 1). The O-Ligs were included into one new CAZy family, GT121 with more than 16,700 members distributed in the four subfamilies.

The greater diversity of the GT121 O-Ligs compared to the GT119 peptidoglycan polymerases and GT120 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<sup>15</sup>. Taxonomically, the GT121 O-Lig family is present in most bacteria, including both Gram-negative and Gram-positive bacteria. 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<sup>20</sup>.

A recently discovered O-Lig, WadA, is bimodular with a GT121 domain appended to a globular glycosyltransferase domain of family GT25, which adds the last sugar to the oligosaccharide core<sup>27</sup>. We have constructed a tree with representative WadA homologs from the GT121 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 GT121 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*.

### Other bacterial polysaccharide polymerases

The fourth functional class of Und-PP-sugar-utilizing GTs are the BP-Pols. As previously mentioned, there is only one experimentally characterized BP-Pol<sup>18</sup>, but several studies have identified BP-Pols from the polysaccharide



**Fig. 1 | Clustering of BP-Pol sequences.** **a** The first step of the clustering: SSN network with nodes representing individual proteins and edges representing pairwise alignment bit scores. Proteins are linked by edges if they have a pairwise score above 110. The resulting clusters are sorted according to number of protein members, with the largest cluster in the upper left corner. **b** The second step of the clustering: HMM models were built for each SSN cluster and the HMMs were compared using HHblits. A network was built with nodes representing SSN clusters and edges representing HHblits scores. SSN clusters are linked by edges if they have

an HHblits score higher than 160. The resulting clusters are referred to as super-clusters and are sorted according to number of SSN clusters. There are two edges between nodes, when the HHblits score is above 160 in both directions. The size of the nodes represents the number of members in the SSN cluster. The 14 largest superclusters (>150 GenBank members) define CAZy families GT122 - GT135. Nodes are colored consistently according to their respective CAZy family in both (a, b).

gene clusters, and we decided to build our families based on these published reports. We thus collected 363 predicted BP-Pol sequences from seven studies for various species, both Gram-negative and Gram-positive bacteria: *Escherichia coli*<sup>28</sup>, *Shigella boydii*, *Shigella dysenteriae*, *Shigella flexneri*<sup>29</sup>, *Salmonella enterica*<sup>30</sup>, *Yersinia pseudotuberculosis*, *Yersinia similis*<sup>31</sup>, *Pseudomonas aeruginosa*<sup>16</sup>, *Acinetobacter baumannii*, *Acinetobacter nosocomialis*<sup>32</sup> and *Streptococcus pneumoniae*<sup>19</sup> (Supplementary Data 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<sup>16</sup>. 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 transmembrane 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 363 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<sup>33</sup>. We then combined the SSN clusters into superclusters in a network analysis based on the HHblits scores (Fig. 1b), resulting in 27 superclusters of varying sizes and 86 singleton clusters. Interestingly, the BP-Pols clustered across taxonomy, and even BP-Pols from Gram-positive and Gram-negative bacteria clustered together. The 14 largest superclusters define new GT families in the CAZy database (GT122-GT135) with a number of members ranging from 159 to 5979 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 taxonomic range, and outside of the taxonomic orders of the original seeds. Several of the families contain members from both Gram-positive and Gram-negative bacteria, for example GT122, GT130, and GT134.

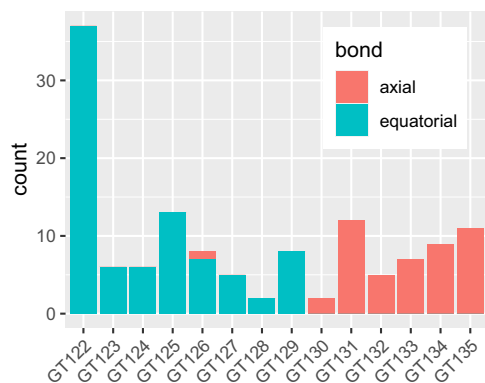
As a way of evaluating our families, we performed structural superimpositions of AlphaFold models of distantly related members of each family. As an example, superimpositions of five distantly related members of GT122 are shown in Supplementary Fig. 3. The sequence identity between these members is relatively low (between 21.4 and 24.3%). Yet, they still produce a meaningful superimposition, and notably, the conserved residues are oriented very similarly.

### 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<sup>16,19,29-32,34</sup>. Additionally, nine sugar structures were included, which were published after the review papers<sup>35-39</sup>. Out of the 150 BP-Pol seed sequences that were included in the new CAZy families, we matched 131 with a sugar structure. The repeat units are oligosaccharides with 3–7 monomers within the backbone, often with branches. In most of the cases, the bond which is formed by the polymerase has been identified in the review papers based on the other GTs in the gene cluster which assemble the repeat units.

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

We found that the stereochemical outcome of BP-Pols appears well conserved within the new BP-Pol CAZy families and varies from one family to another (Fig. 2). There is only one exception; in family GT126, the polymerase linkages are all equatorial except for the O-antigen in



**Fig. 2 | Level of conservation of stereochemical outcome of the reaction catalyzed in the various BP-Pol families.** The bars represent the number of enzymes that are known to employ either retaining (making an axial bond) or inverting (making an equatorial bond) mechanisms in each of the new BP-Pol families.

*Pseudomonas aeruginosa* O4, where it is axial. This could be due to an error in the chemical structure or the serotype designation or that the *P. aeruginosa* O4 polymerase constitutes an exception.

Next, we investigated whether there was a correlation between the structures of the transferred sugars and the sequence similarity of the BP-Pols. We created phylogenetic trees of the BP-Pols in each family and visualized them with the corresponding transferred repeat units. We observe that the sugars within each family show similarity and this similarity appears to correlate with the structure of the tree, implying that polymerases with similar sequence utilize similar substrates (Fig. 3, Supplementary Fig. 4). The ends of the repeat units, i.e. the subsite moieties immediately upstream (+1) and downstream (-1) of the newly created bond (Fig. 4) seem to be most conserved whereas more variability occurs in the middle part. We hypothesize that the +1 and -1 subsites are the moieties most important for recognition by the active site of the BP-Pol.

We observe examples of BP-Pols from distant taxonomic origin that cluster in the same CAZy family and have highly similar sugars. For example, *Escherichia coli* O178 and *Streptococcus pneumoniae* 47A in GT125 transfer sugars with almost identical backbones, suggestive of horizontal gene transfer (Fig. 3). There is only a slight variance in the middle of the repeat unit. This suggests that there is less constraints on the central part of the repeat unit than on the extremities that define the donor and the acceptor.

We next attempted to quantify the correlation between BP-Pol sequence and carbohydrate structure. For this we developed an original pairwise oligosaccharide similarity score. In our scoring scheme, the similarity of two glycans is estimated by examining the -1 and +1 subsites, as we expect that these are the moieties most fitting the active site of the BP-Pol (Fig. 4). The 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).

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

### Comparison of families

Others have previously reported sequence and structural similarity between SEDS, O-Lig and some BP-Pols<sup>13,14,21,23</sup>. In order to investigate the relatedness of the new CAZy families, we compared the family HMMs by all-vs-all HHblits analyses<sup>33</sup> (Fig. 6). Strikingly, we observe that the retaining BP-Pol families cluster together on the heatmap along with the retaining ECA-Pols,

while the inverting BP-Pols form two distinct groups, one of them containing the inverting SEDS (GT119) and the inverting O-Ligs (GT121). The background noise between some inverting and retaining enzymes is likely due to the general conservation of the successive transmembrane helices, which is altered in the GT122-GT123-GT124 subgroup due to their different architecture (see below).

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<sup>40</sup>. In extension of the report of the GT-C<sub>B</sub> class by Alexander and Locher<sup>23</sup>, and based on the above-mentioned similarities between the new CAZy families, we can now define three sequence-based clans: GT-C<sub>1</sub> consisting of inverting BP-Pol families, SEDS and O-Lig, GT-C<sub>2</sub> consisting of retaining BP-Pol families and ECA-Pol, and GT-C<sub>3</sub> consisting of inverting BP-Pol families (Table 1, Fig. 6). 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 on the sequence similarity to SEDS and O-Ligs, we have assigned the BP-Pol families of clan GT-C<sub>1</sub> to the structural subclass GT-C<sub>B</sub> of Alexander and Locher<sup>23</sup>. In addition, we also present in Table 1 the families of GT-C glycosyltransferases that have not yet been assigned to a structural class.

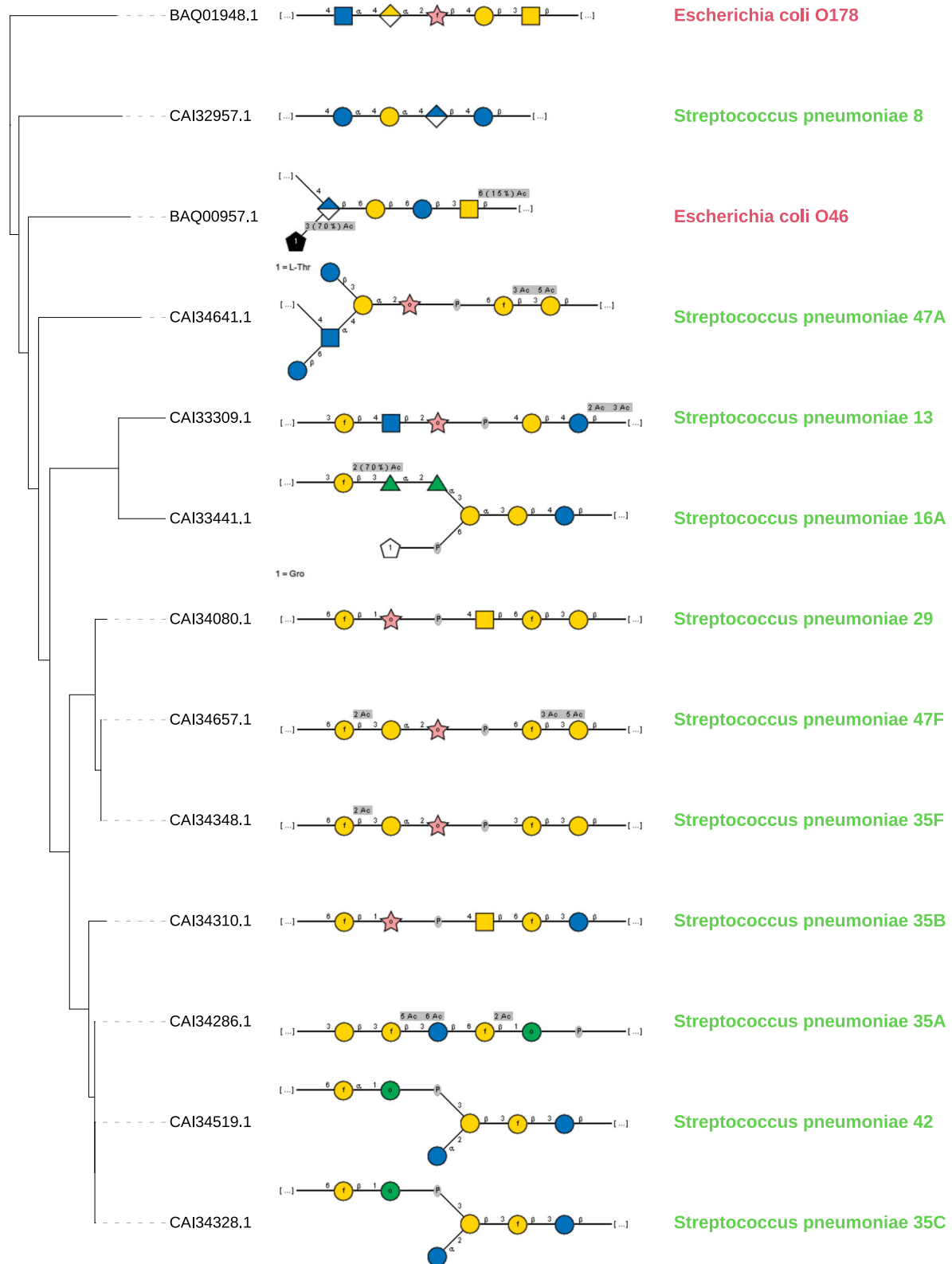
We then examined residue conservation and the general architecture of the enzymes in the clans. Based on the above mentioned pairwise HHblits analyses and structural superimpositions (Supplementary Figs. 5–7), we tried to evaluate which architectural features and conserved residues are common within the clans. Indeed, there are some common features across most families. In all the families, all the conserved residues are located on the outer face of the membrane (Fig. 7). Enzymes of clans GT-C<sub>1</sub> and GT-C<sub>2</sub> have a long extracellular loop close to the C-terminus containing conserved residues (Fig. 7). In stark contrast, families GT122, GT123 and GT124 of clan GT-C<sub>3</sub> have an architecture completely different from that of the two other clans (Fig. 7), with a long loop located close to the N-terminus.

The families in GT-C<sub>1</sub> show a distinct pattern of residue conservation. 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<sup>21</sup>. The SEDS family (GT119) also has a conserved Arg which aligns with the second conserved Arg in O-Lig and a conserved essential Asp which aligns with the conserved His in O-Lig (Fig. 7). Likewise, all the BP-Pols in the clan have 1–2 conserved Args, some of which align to the O-Lig Args in the HHblits alignments, and we hypothesize that they also play the role of binding to the diphosphate. Similarly, all the families in the clan except for GT127 have either a conserved Asp or Glu, which align with the conserved His of O-Lig and the conserved Asp of SEDS (Fig. 7). We hypothesize that these Glu and Asp residues also play the role of activating the acceptor. As an example, the superimposition of the published O-Lig structure (7TPG)<sup>21</sup> and an AlphaFold model from one representative of the inverting BP-Pol family GT126 is shown in Fig. 8a. 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 and Glu are in the same position. As mentioned above, GT127 does not have a conserved Asp or Glu in the same position as the rest of the families. However, it has a conserved Glu in a loop between transmembrane helices 5 and 6, which likely plays the same role.

In the retaining clan GT-C<sub>2</sub>, the pattern of conservation is different. Here, most of the families have 2–3 conserved Arg/Lys and 1–2 conserved Tyr (Fig. 7). Interestingly, we observe that the ECA-Pol family GT120 shows high similarity with one of the BP-Pol families, GT134. A superimposition of AlphaFold models from each family shows that the conserved residues are oriented very similarly, despite the low overall similarity (RMSD 5.4 Å over 360 residues) (Fig. 8b).

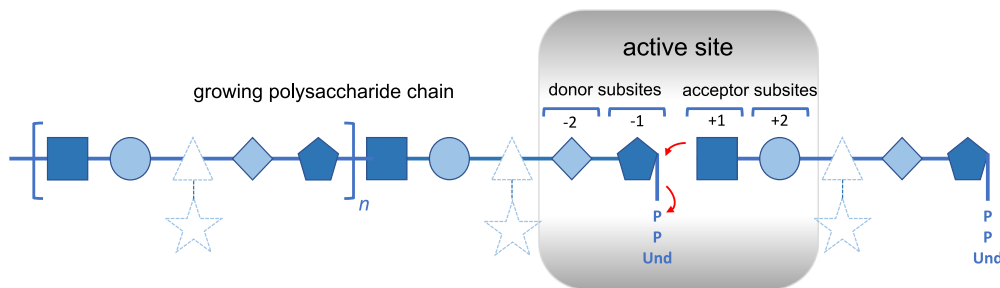
The families in the inverting clan GT-C<sub>3</sub> all have two conserved Arg, a conserved Asp, and a conserved His, all of which align between the families in the HHblits alignments (Fig. 7).



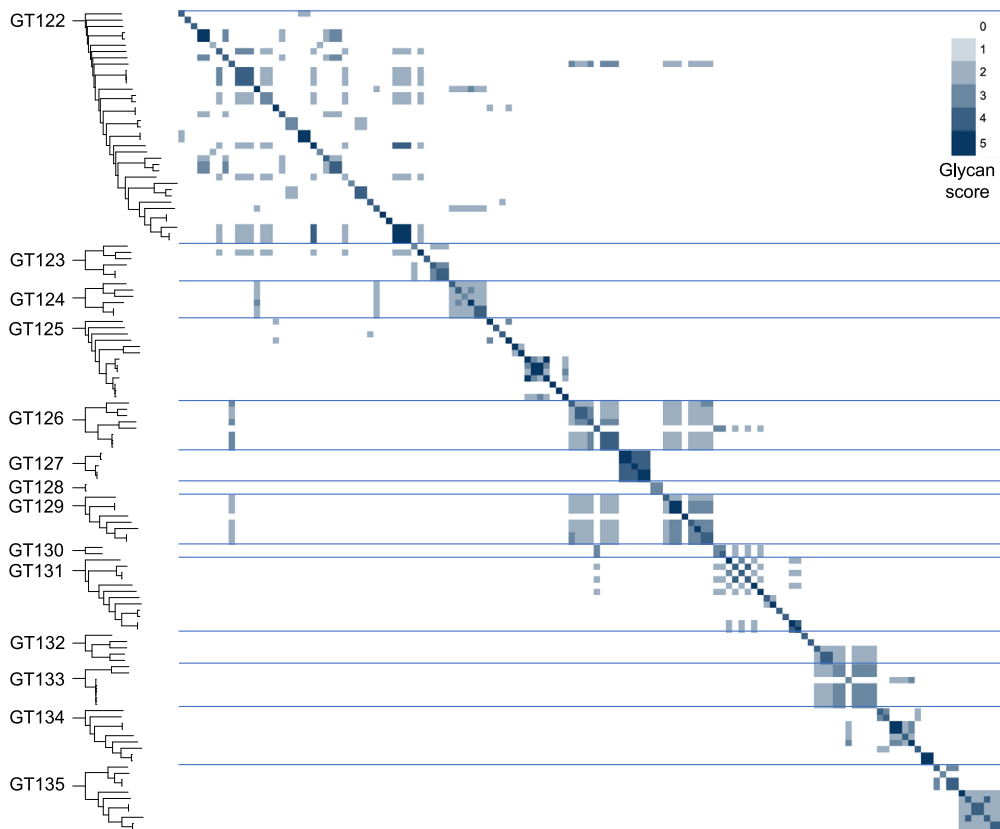


**Fig. 3 | Comparison of repeat unit sugars transferred by BP-Pols in GT125.** The transferred repeat unit structures (in SNFG representation) are shown on a phylogenetic tree of BP-Pols in family GT125. There is an overall similarity between all the transferred sugars in the family and the similarity appears to correlate with the tree structure, i.e., BP-Pol sequence 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 where the enzyme does not interact with the sugar. Note that the +1 site corresponds to the non-reducing end of the depicted sugar structures and the -1 site corresponds to the reducing end. Notably, the family contains BP-Pols from distant taxonomic origin and that yet transfer similar repeat units.



**Fig. 4 | An idealized representation of a BP-Pol.** The donor is the growing glycan chain activated by Und-PP while the acceptor is a single repeat unit linked to Und-PP. 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.



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

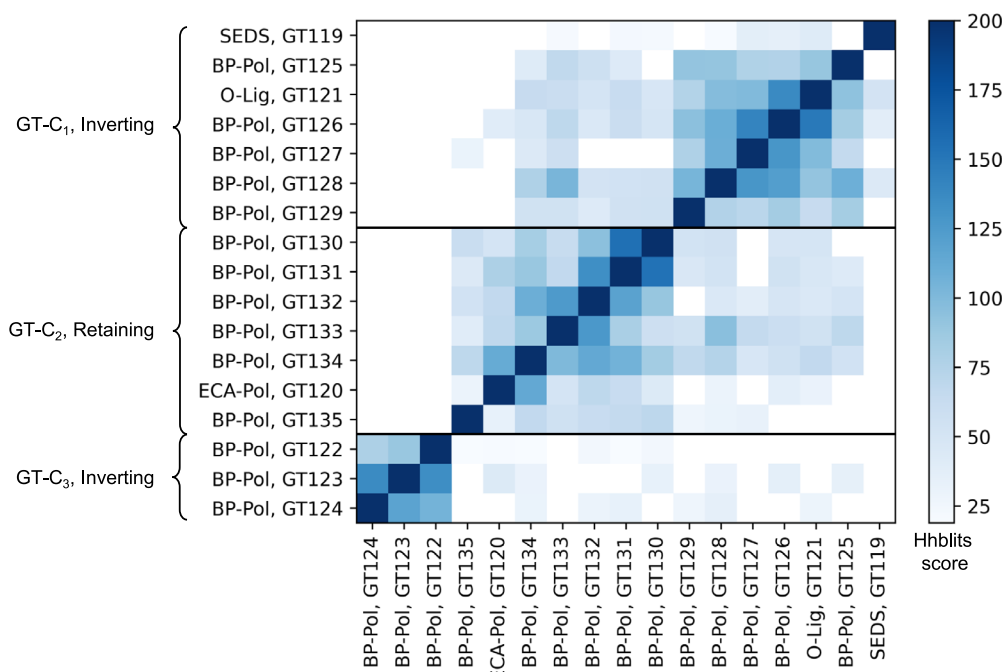
## Discussion

Here we have added 17 glycosyltransferase families (GT119 to GT135) to the CAZy database bringing the total of covered families from 118 to 135. 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.

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 previously because of their extreme diversity even within strains of a single species<sup>28</sup>, 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 current sequence diversity, we were able to form larger families of similar polymerases from widely different taxonomies, thereby revealing conserved residues that are most likely functionally important.

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* O108 (GenBank BAQ01516.1) and *A. baumannii* O24 (GenBank AHB32586.1)



**Fig. 6 | Relatedness of the new CAZy families and definition of clans.** Inter-family HHblits bit scores are shown in a heatmap 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. The

inverting BP-Pols form two clans, GT-C<sub>1</sub> which also contains the inverting SEDS (GT119) and the inverting O-Ligs (GT121) and GT-C<sub>3</sub> containing only BP-Pols. The retaining BP-Pol families form one clan, GT-C<sub>2</sub>, which also contains the retaining ECA-Pol family (GT120).

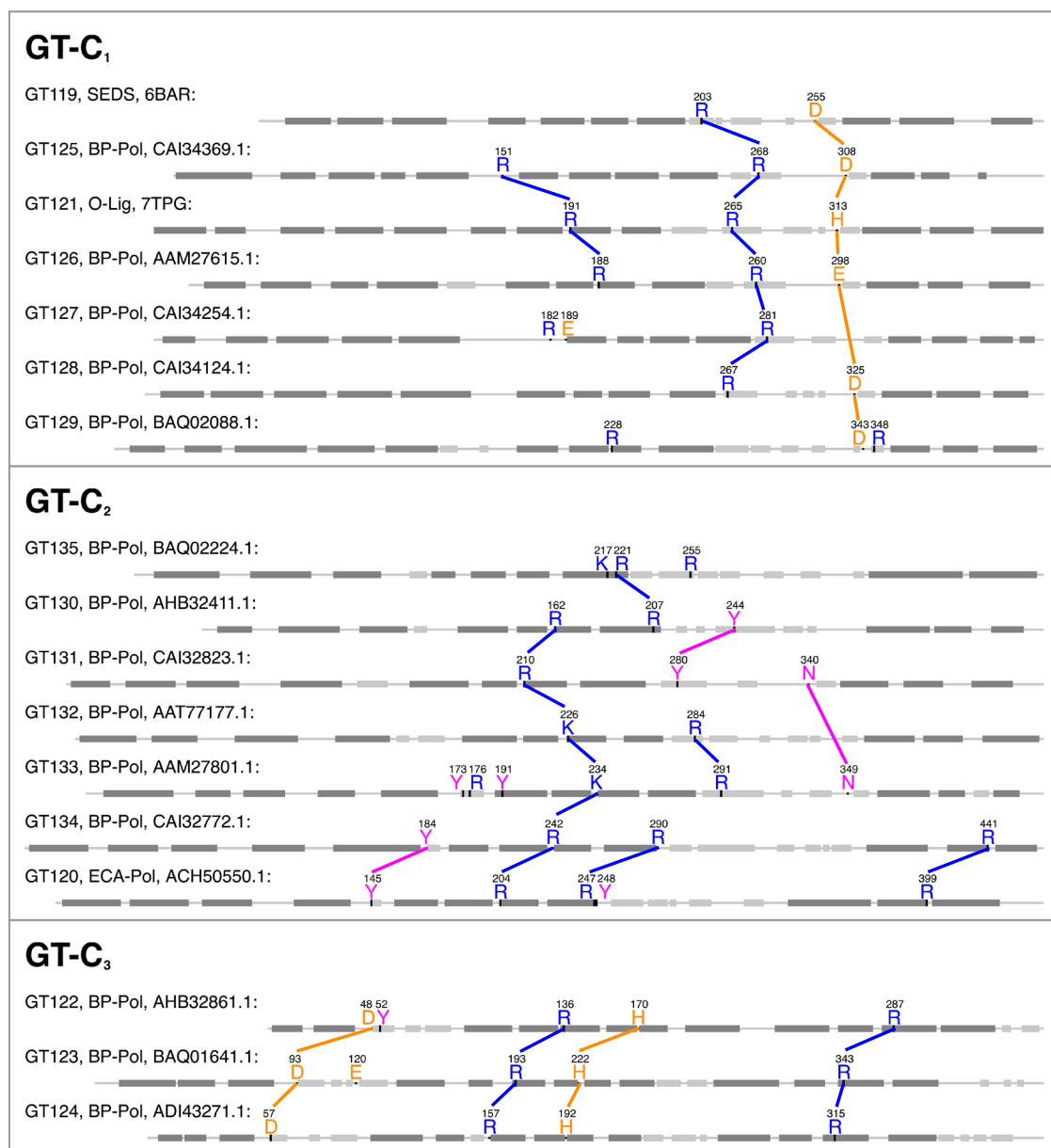
**Table 1 | Structural subclasses, clans and families of GT-C fold glycosyltransferases and relationships to mechanism and glycosyl donor**

Structural subclass Alexander & Locher	CAZy clan	CAZy families	Mechanism	Donor
GT-C <sub>A</sub>	-	GT53	Inverting	Lipid-P-monosaccharide
GT-C <sub>A</sub>	-	GT83	Inverting	Lipid-P-monosaccharide
GT-C <sub>A</sub>	-	GT39	Inverting	Lipid-P-monosaccharide
GT-C <sub>A</sub>	-	GT57	Inverting	Lipid-P-monosaccharide
GT-C <sub>A</sub>	-	GT66	Inverting	Lipid-PP-oligosaccharide
GT-C <sub>B</sub>	GT-C <sub>1</sub>	GT119, GT121, GT125, GT126, GT127, GT128, GT129	Inverting	Lipid-PP-oligosaccharide
-	GT-C <sub>2</sub>	GT120, GT130, GT131, GT132, GT133, GT134, GT135	Retaining	Lipid-PP-oligosaccharide
-	GT-C <sub>3</sub>	GT122, GT123, GT124	Inverting	Lipid-PP-oligosaccharide
-	-	GT22	Inverting	Lipid-P-monosaccharide
-	-	GT50	Inverting	Lipid-P-monosaccharide
-	-	GT58	Inverting	Lipid-P-monosaccharide
-	-	GT59	Inverting	Lipid-P-monosaccharide

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.

Moreover, we observe that the sequence diversity within the families we have built is minimal for peptidoglycan polymerases (GT119) and ECA-Pols (GT120), and then increases gradually for O-Ligs (GT121) and is maximal for BP-Pols (GT122-GT135). We hypothesize that sequence diversity reflects the donor and acceptor diversity in each family since the latter increases accordingly; the enzymes in the SEDS and ECA-Pol families act with the same donor and same acceptor, the enzymes in the O-Lig family act with different donors but same acceptor, and for the enzymes in the BP-Pol families act on different donors and different acceptors.

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 (S<sub>N</sub>2) or below (S<sub>N</sub>i) the sugar ring of the donor<sup>4</sup>. 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<sup>41</sup>. 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



**Fig. 7 | Equivalent conserved residues in the clans.** Conserved residues of each of the new CAZy families are shown on sequences of representative family members. Colored lines are shown between conserved residues from different families, which align in HHblits alignments and co-localize in structural superimpositions (Supplementary Figs. 5–7). Transmembrane helices are shown in dark gray boxes,

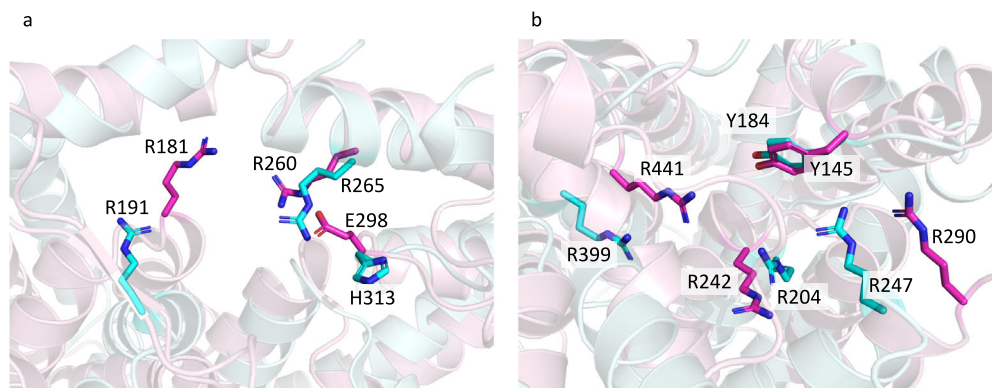
extracellular helices are shown in light gray boxes. The secondary structures were taken from the crystal structures for family GT119 and GT121 (6BAR and 7TPG respectively) and from AlphaFold models for all other families. The R210 in GT131 is either K or R in the family.

suggests that retaining BP-Pols also operate by an  $S_Ni$  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 GT120 and GT130–GT135. Additionally, the  $S_Ni$  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<sup>23</sup>. The GT families that we describe here significantly expand the GT-C class in the CAZy database ([www.cazy.org](http://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<sup>4</sup>. 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 GT122, GT123, and GT124 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, GT122, GT123 and GT124 were assigned to clan GT-C<sub>3</sub>.





**Fig. 8 | Structural superimpositions of members of different functional classes belonging to the same clans.** **a** Superimposition of O-Lig in cyan (GT121, PDB: 7TPG) and AlphaFold model of BP-Pol in pink (GT126, Genbank accession: AAM27615.1) showing that the conserved Glu in the BP-Pol aligns with the conserved His in the O-Lig, which has been proposed to activate the acceptor<sup>21</sup> (RMSD

5.3 Å over 192 residues, sequence identity 20.8% over 485 residues). **b** Structural superimpositions of AlphaFold models of ECA-Pol in cyan (GT120, Genbank accession: ACH50550.1) and BP-Pol in pink (GT134, Genbank accession: CAI32772.1) illustrating structural similarity and co-localization of the conserved residues (RMSD 5.4 Å over 360 residues, sequence identity 17.1% over 543 residue).

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 a certain 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 Und-PP of the donor, i.e. the residues closest to the reaction center (Fig. 4). 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 the presence of a phosphodiester bond.

The version of the glycan similarity score presented here was inspired in part by observed structural similarities in different O-antigen repeat units assembled by very similar BP-Pols<sup>16</sup>. The repeat unit comparison involves a translation of glycan IUPAC nomenclature to a reduced alphabet of terms representing only backbone configuration, i.e., ignoring chemical modifications and sidechains. Furthermore, a positive similarity score requires an entire identical match of all backbone elements at both donor and acceptor positions (−1 and +1 sites in Fig. 4, respectively). Despite these simplifications, the similarity score reveals, with exceptions, an overall greater intra- rather than inter-family oligosaccharide similarity (Fig. 5). 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<sup>42</sup>. The multiplicity of polysaccharide biosynthesis loci in some genomes makes it sometimes difficult to assign a particular polysaccharide structure to a particular biosynthesis operon.

We observed that the O-Lig family (GT121) was present in many Gram-positive bacteria such as *Streptococcus pneumoniae*. 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<sup>17,43</sup>. Thus a hypothesis could be that the GT121 members in *S. pneumoniae* are responsible for the ligation of CPS to the peptidoglycan layer in these bacteria.

As already shown in other occasions, the sequence-based classification of carbohydrate-active enzymes of the CAZy database has predictive power. The case of the GT families described here supports this view as the invariant

residues in the families not only co-localize in the same area of the three-dimensional structures (whether actual or AlphaFold-predicted), but also correspond to the residues found essential for function in the families where this has been studied experimentally. The families described herein also show mechanistic conservation and thus the stereochemistry of glycosyl transfer can be predicted. Finally, the observed similarity in oligosaccharide repeat units that accompanies sequence similarity has also predictive power and paves the way to the future possibility of in silico serotyping based on DNA sequence.

## Methods

### Alignment-based Clustering (Aclust)

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

### Building the peptidoglycan polymerase family (GT119)

The peptidoglycan polymerase family, GT119, was built by using Blastp from BLAST+ 2.12.0+<sup>48</sup> with the sequences of the characterized SEDS proteins (PDB 6BAR, 8TJ3, 8BH1 and GenBank accession CAB15838.1) against GenBank with a threshold of approximately 30% to retrieve the family members. Next, an MSA was generated with MAFFT v7.508 using the L-INS-i strategy<sup>45</sup>, and an HMM model was built with hmmbuild of HMMER 3.3.2<sup>49</sup>. The family was further populated using hmmsearch from HMMER 3.2.2 against GenBank.

### Building the enterobacterial common antigen polymerase family (GT120)

A sequence library of ECA-Pols was constructed by using Blastp with the seed sequence (GenBank accession AAC76800.1) against the NCBI non-

redundant database version 61 with an E-value threshold of  $1e-60$ . The hits were redundancy reduced using CD-HIT 4.8.1<sup>50</sup> with a threshold of 99%. The redundancy-reduced pool of ECA-Pol sequences was clustered using our in-house tool Aclust (see above), and the tree showed one large clade and a few outliers. All the sequences in the large clade were used to build an MSA using MAFFT v7.508 with the L-INS-i strategy<sup>45</sup>. An HMM was built based on this MSA using hmmbuild of HMMER 3.3.2<sup>49</sup>. The family GT121 was built in CAZy and populated using Blastp against GenBank with an approximate threshold of 30% and hmmsearch against GenBank.

### Building the O-antigen ligase family (GT121)

37 O-Lig sequences were selected from literature (Supplementary Data 1) and expanded using Blastp against the NCBI non-redundant database with an E-value cut-off of  $1e-60$ . Redundancy reduction was performed on the resulting sequence pool using CD-HIT with a threshold of 99%, resulting in a pool of 1402 sequences. A phylogenetic tree of the pool of O-Lig sequences was generated using Aclust (see above), which showed deep clefts between main branches, and branches with sufficient internal diversity (Supplementary Fig. 2). Based on these results, four subfamilies were determined. An MSA was built for the family as well as for the subfamilies with MAFFT v7.508 using the L-INS-i strategy. HMMs were built based on the MSAs using the hmmbuild of HMMER 3.3.2<sup>49</sup>. The family was populated using Blastp against GenBank with an approximate threshold of 30% identity with the seed sequences and using hmmsearch with the family and subfamily HMMs.

### Building the bacterial polysaccharide polymerase families (GT122-GT135)

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

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

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

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

### Analysis of sugar repeat unit structures

In order to analyze the relation between BP-Pol sequence and structure of the transferred repeat unit, we retrieved the repeat unit structures for the serotypes for the BP-Pols that were included in the new CAZy families. The repeat unit structures were retrieved from the same review papers from which we retrieved the BP-Pol sequences<sup>16,19,29-32</sup>, except for the sugars for *E. coli*, where the sugar structures have been reported elsewhere<sup>34</sup>. Nine additional repeat unit structures were included for *S. pneumoniae*, which were published after the review paper; serotypes 16A<sup>35</sup>, 33A<sup>36</sup>, 33C and 33D<sup>37</sup>, 35C and 35F<sup>38</sup>, 42 and 47F<sup>56</sup> and 47A<sup>57</sup>. For *Y. pseudotuberculosis* O3 and *S. pneumoniae* 33B, we used the revised structures<sup>37,39</sup>. *Pseudomonas aeruginosa* O2 and O16 contain two BP-Pol genes; one BP-Pol localized in the O-antigen biosynthesis cluster, which polymerizes the sugar repeat units with an  $\alpha$  bond and one BP-Pol localized outside the biosynthesis cluster which polymerizes the repeat units with a  $\beta$  bond<sup>58</sup>. Since the BP-Pols reported in<sup>16</sup> are from the O-antigen cluster, we report the sugar structure with the  $\alpha$  bond.

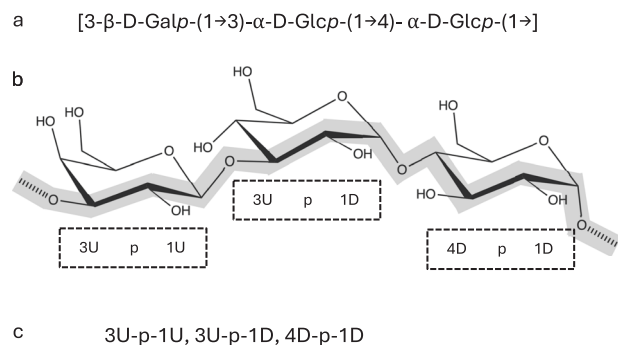
The linkages formed by the polymerase have been determined in all of these papers, except for a few cases. This determination is based on the other GTs in the gene cluster, in particular the initial GT which transfers the first monosaccharides to the Und-PP anchor. The cases where the polymerase linkage has not been unambiguously determined in the review papers are *E. coli* O166, O78, O152, O81, O83, O11, O112ab, O167, O187, O142, O117, O107, O185, O42, O28ac, O28ab, for which there are two or more possible polymerase linkages. For the structures that were published after the review papers, the polymerase bond had not been determined in *S. pneumoniae* 33A and 47A. For *S. pneumoniae* 33A, we determined the linkage based on the presence of the initial transferase *wchA* in the gene cluster, which transfers a glucose-1-phosphate to Und-PP<sup>19</sup>. In *S. pneumoniae* 47A the initial transferase is WcjG, which transfers Galp or Galf<sup>19</sup>. Since the repeat unit contains both Gal and Galp, we could not determine the polymerase linkage unambiguously. However, the repeat unit is very similar to other repeat units in the family (most similar to that of *S. pneumoniae* 13), and we proposed the equivalent polymerase linkage.

The CSDB database (<http://csdb.glycoscience.ru>)<sup>59</sup> was used to retrieve literature, SNFG image representations and linear sugar strings of the repeat unit structures. Phylogenetic trees for BP-Pol families with sugar structures were generated using MAFFT v7.508<sup>45</sup> with the L-INS-i strategy to supply an initial multiple sequence alignment, followed by Aclust (see above) for distance matrix embedding and clustering. The trees were visualized in iTOL<sup>60</sup>. The barplot was generated using R<sup>61</sup>, Rstudio<sup>62</sup>, and the ggplot2 package<sup>63</sup>.

### 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 (Fig. 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 direction, adding one to the score for each additional identical match, but terminating at the first non-identity or possible re-use of a backbone position.

To facilitate comparison, oligosaccharide sequences are translated from IUPAC nomenclature into symbols that represent elements of backbone geometry, only considering monomer dimension and stereochemistry of acceptor and anomeric donor carbon atoms, and ignoring sidechains and chemical modification (Fig. 9). Briefly, the monomer dimension is represented by a single letter P, F or L depending on whether the monomer sugar is a pyranose, furanose or is linear, respectively. Stereochemistry of the acceptor and donor carbon atoms is represented by the index number of the carbon position within the ring/monomer, followed by a single letter U, D or N depending on whether the linked oxygen atom is U (up=above the monomer ring),



**Fig. 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 gray line) and transfer bond (hatched line segments), and translated geometric subunits below. c Completed translation.**

D (down=below the monomer ring), or N (neither above or below the ring). The N symbol is assigned in cases of conformational flexibility such as with alditols or C6 linkages. At present, in scoring the similarity of two thus translated residues, the entirety of the translation strings must be identical to achieve a score of +1. Further details and limitations will be presented elsewhere (G.P. Gippert, manuscript in preparation).

### Comparison of the families

Pairwise HHblits analyses<sup>33</sup> were performed for each of the new CAZy families. The HHblits scores were visualized in a heatmap using Python Matplotlib<sup>64</sup>.

AlphaFold2<sup>14</sup> structures were generated of representative proteins from the families using the ColabFold implementation<sup>65</sup> on our internal GPU cluster processed with the recommended settings. The best ranked relaxed model was used. The protein structures were visualized in PyMOL<sup>66</sup> and pairwise structural superimpositions were performed using the CEalign algorithm<sup>67</sup>.

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

### Data availability

Accessions to the seed sequences utilized in this work are given in Supplementary Data 1 and 2; the constantly updated content of families GT119 - GT135 is given in the online CAZy database at [www.cazy.org](http://www.cazy.org).

### Code availability

Source code for Aclust may be obtained via GitHub at <https://github.com/GarryGippert/Aclust>.

Received: 19 September 2023; Accepted: 16 February 2024;

Published online: 07 March 2024

### References

- Varki, A. et al. (eds.) *Essentials of Glycobiology* (Cold Spring Harbor Laboratory Press, Cold Spring Harbor (NY), 2022), 4th edn.
- Laine, R. A. A calculation of all possible oligosaccharide isomers both branched and linear yields  $1.05 \times 10^{12}$  structures for a reducing hexasaccharide: the Isomer Barrier to development of single-method saccharide sequencing or synthesis systems. *Glycobiology* **4**, 759–767 (1994).
- Lapébie, P., Lombard, V., Drula, E., Terrapon, N. & Henrissat, B. Bacteroidetes use thousands of enzyme combinations to break down glycans. *Nat. Commun.* **10**, 2043 (2019).
- Lairson, L., Henrissat, B., Davies, G. & Withers, S. Glycosyltransferases: structures, functions, and mechanisms. *Annu. Rev. Biochem.* **77**, 521–555 (2008).
- Campbell, J. A., Davies, G. J., Bulone, V. & Henrissat, B. A classification of nucleotide-diphospho-sugar glycosyltransferases based on amino acid sequence similarities. *Biochem. J.* **326**, 929–939 (1997).
- Drula, E. et al. The carbohydrate-active enzyme database: functions and literature. *Nucleic Acids Res.* **50**, D571–D577 (2022).
- McDonald, A. G. & Tipton, K. F. Fifty-five years of enzyme classification: advances and difficulties. *FEBS J.* **281**, 583–592 (2014).
- Coutinho, P. M., Deleury, E., Davies, G. J. & Henrissat, B. An evolving hierarchical family classification for glycosyltransferases. *J. Mol. Biol.* **328**, 307–317 (2003).
- Knauer, R. & Lehle, L. The oligosaccharyltransferase complex from yeast. *Biochim. et Biophys. Acta* **1426**, 259–273 (1999).
- Cho, H. Assembly of bacterial surface glycopolymers as an antibiotic target. *J. Microbiol.* **60**, 359–367 (2023).
- Sjodt, M. et al. Structure of the peptidoglycan polymerase RodA resolved by evolutionary coupling analysis. *Nature* **556**, 118–121 (2018).
- Käshammer, L. et al. Cryo-EM structure of the bacterial divisome core complex and antibiotic target FtsW/IBL. *Nat. Microbiol.* **8**, 1149–1159 (2023).
- Nygaard, R. et al. Structural basis of peptidoglycan synthesis by *E. coli* RodA-PBP2 complex. *Nat. Commun.* **14**, 5151 (2023).
- Meeske, A. J. et al. SEDS proteins are a widespread family of bacterial cell wall polymerases. *Nature* **537**, 634–638 (2016).
- Di Lorenzo, F. et al. A journey from structure to function of bacterial lipopolysaccharides. *Chem. Rev.* **122**, 15767–15821 (2022).
- Islam, S. T. & Lam, J. S. Synthesis of bacterial polysaccharides via the Wzx/Wzy-dependent pathway. *Can. J. Microbiol.* **60**, 697–716 (2014).
- Whitfield, C., Wear, S. S. & Sande, C. Assembly of bacterial capsular polysaccharides and exopolysaccharides. *Annu. Rev. Microbiol.* **74**, 521–543 (2020).
- Woodward, R. et al. In vitro bacterial polysaccharide biosynthesis: defining the functions of Wzy and Wzz. *Nat. Chem. Biol.* **6**, 418–423 (2010).
- Bentley, S. D. et al. Genetic analysis of the capsular biosynthetic locus from All 90 pneumococcal serotypes. *PLoS Genet.* **2**, e31 (2006).
- 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).
- Ashraf, K. U. et al. Structural basis of lipopolysaccharide maturation by the O-antigen ligase. *Nature* **604**, 371–376 (2022).
- Rai, A. K. & Mitchell, A. M. Enterobacterial common antigen: synthesis and function of an enigmatic molecule. *mBio* **11**, 1–19 (2020).
- Alexander, J. A. N. & Locher, K. P. Emerging structural insights into C-type glycosyltransferases. *Curr. Opin. Struct. Biol.* **79**, 102547 (2023).
- Emami, K. et al. RodA as the missing glycosyltransferase in *Bacillus subtilis* and antibiotic discovery for the peptidoglycan polymerase pathway. *Nat. Microbiol.* **2**, 16253 (2017).
- Maczuga, N., Tran, E. N. H., Qin, J. & Morona, R. Interdependence of *Shigella flexneri* O Antigen and enterobacterial common antigen biosynthetic pathways. *J. Bacteriol.* **204**, e00546–21 (2022).
- Lombard, V., Golaconda Ramulu, H., Drula, E., Coutinho, P. M. & Henrissat, B. The carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res.* **42**, D490–D495 (2014).



27. Servais, C. et al. Lipopolysaccharide biosynthesis and traffic in the envelope of the pathogen *Brucella abortus*. *Nat. Commun.* **14**, 911 (2023).
28. Iguchi, A. et al. A complete view of the genetic diversity of the *Escherichia coli* O-antigen biosynthesis gene cluster. *DNA Res.* **22**, 101–107 (2015).
29. Liu, B. et al. Structure and genetics of *Shigella* O antigens. *FEMS Microbiol. Rev.* **32**, 627–653 (2008).
30. Liu, B. et al. Structural diversity in *Salmonella* O antigens and its genetic basis. *FEMS Microbiol. Rev.* **38**, 56–89 (2014).
31. Kenyon, J. J., Cunneen, M. M. & Reeves, P. R. Genetics and evolution of *Yersinia pseudotuberculosis* O-specific polysaccharides: a novel pattern of O-antigen diversity. *FEMS Microbiol. Rev.* **41**, 200–217 (2017).
32. Hu, D., Liu, B., Dijkshoorn, L., Wang, L. & Reeves, P. R. Diversity in the Major Polysaccharide Antigen of *Acinetobacter Baumannii* Assessed by DNA Sequencing, and Development of a Molecular Serotyping Scheme. *PLoS ONE* **8**, e70329 (2013).
33. Remmert, M., Biegert, A., Hauser, A. & Söding, J. HHblits: lightning-fast iterative protein sequence searching by HMM-HMM alignment. *Nat. Methods* **9**, 173–175 (2012).
34. Liu, B. et al. Structure and genetics of *Escherichia coli* O antigens. *FEMS Microbiol. Rev.* **44**, 655–683 (2020).
35. Li, C. et al. Structural, biosynthetic, and serological cross-reactive elucidation of capsular polysaccharides from *Streptococcus pneumoniae* Serogroup 16. *J. Bacteriol.* **201**, 13 (2019).
36. Lin, F. L. et al. Identification of the common antigenic determinant shared by *Streptococcus pneumoniae* serotypes 33A, 35A, and 20 capsular polysaccharides. *Carbohydr. Res.* **380**, 101–107 (2013).
37. Lin, F. L. et al. Structure elucidation of capsular polysaccharides from *Streptococcus pneumoniae* serotype 33C, 33D, and revised structure of serotype 33B. *Carbohydr. Res.* **383**, 97–104 (2014).
38. Bush, C. A., Cisar, J. O. & Yang, J. Structures of capsular polysaccharide serotypes 35F and 35C of *Streptococcus pneumoniae* determined by nuclear magnetic resonance and their relation to other cross-reactive serotypes. *J. Bacteriol.* **197**, 2762–2769 (2015).
39. Kondakova, A. N. et al. Reinvestigation of the O-antigens of *Yersinia pseudotuberculosis*: revision of the O<sub>2c</sub> and confirmation of the O<sub>3</sub> antigen structures. *Carbohydr. Res.* **343**, 2486–2488 (2008).
40. Henrissat, B. & Bairoch, A. Updating the sequence-based classification of glycosyl hydrolases. *Biochem. J.* **316**, 695–696 (1996).
41. Doyle, L. et al. Mechanism and linkage specificities of the dual retaining  $\beta$ -Kdo glycosyltransferase modules of KpsC from bacterial capsule biosynthesis. *J. Biol. Chem.* **299**, 104609 (2023).
42. Krinos, C. M. et al. Extensive surface diversity of a commensal microorganism by multiple DNA inversions. *Nature* **414**, 555–558 (2001).
43. Paton, J. C. & Trappetti, C. *Streptococcus pneumoniae* Capsular Polysaccharide. *Microbiol. Spectr.* **7**, 7.2.33 (2019).
44. Smith, T. & Waterman, M. Identification of common molecular subsequences. *J. Mol. Biol.* **147**, 195–197 (1981).
45. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software Version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
46. Sonnhammer, E. L. & Hollich, V. Scoredist: a simple and robust protein sequence distance estimator. *BMC Bioinform.* **6**, 108 (2005).
47. Crippen, G. & Havel, T. *Distance Geometry and Molecular Conformation*. Chemometrics research studies series (Research Studies Press, 1988).
48. Camacho, C. et al. BLAST+: architecture and applications. *BMC Bioinform.* **10**, 421 (2009).
49. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
50. Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
51. Huszczyński, S. M., Hao, Y., Lam, J. S. & Khursigara, C. M. Identification of the *Pseudomonas aeruginosa* O17 and O15 O-Specific Antigen Biosynthesis Loci Reveals an ABC Transporter-Dependent Synthesis Pathway and Mechanisms of Genetic Diversity. *J. Bacteriol.* **202** <https://doi.org/10.1128/jb.00347-20> (2020).
52. Hagberg, A., Swart, P. & S Chult, D. Exploring network structure, dynamics, and function using networkx. In *Proc. 7th Annual Python in Science Conference, Pasadena, CA, August 19–24, 2008*. 11–16 (2008).
53. Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**, 1189–1191 (2009).
54. Steinegger, M. et al. HH-suite3 for fast remote homology detection and deep protein annotation. *BMC Bioinform.* **20**, 473 (2019).
55. Shannon, P. et al. Cytoscape: a Software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
56. Petersen, B. O., Meier, S., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Determination of native capsular polysaccharide structures of *Streptococcus pneumoniae* serotypes 39, 42, and 47F and comparison to genetically or serologically related strains. *Carbohydr. Res.* **395**, 38–46 (2014).
57. Petersen, B. O., Hindsgaul, O., Paulsen, B. S., Redondo, A. R. & Skovsted, I. C. Structural elucidation of the capsular polysaccharide from *Streptococcus pneumoniae* serotype 47A by NMR spectroscopy. *Carbohydr. Res.* **386**, 62–67 (2014).
58. Lam, J. S., Taylor, V. L., Islam, S. T., Hao, Y. & Kocincová, D. Genetic and functional diversity of *Pseudomonas aeruginosa* lipopolysaccharide. *Front. Microbiol.* **2**, 118 (2011).
59. Toukach, P. V. & Egorova, K. S. Carbohydrate structure database merged from bacterial, archaeal, plant and fungal parts. *Nucleic Acids Res.* **44**, D1229–D1236 (2016).
60. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
61. R Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, Vienna, Austria, 2023).
62. Posit team. *RStudio: Integrated Development Environment for R* (Posit Software, PBC, Boston, MA, 2023).
63. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag New York, 2016).
64. Hunter, J. D. Matplotlib: A 2D graphics environment. *Comput. Sci. Eng.* **9**, 90–95 (2007).
65. Mirdita, M. et al. ColabFold: making protein folding accessible to all. *Nat. Methods* **19**, 679–682 (2022).
66. Schrödinger, L.C.C. The PyMOL Molecular Graphics System, Version 2.5 (2020).
67. Shindyalov, I. N. & Bourne, P. E. Protein structure alignment by incremental combinatorial extension (CE) of the optimal path. *Protein Eng.* **11**, 739–747 (1998).

## Acknowledgements

This work was supported by 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.

### 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.

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s42003-024-05930-2>.

**Correspondence** and requests for materials should be addressed to Bernard Henrissat.

**Peer review information** *Communications Biology* thanks Srijak Bhatnagar and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Primary Handling Editor: Tobias Goris.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2024