**The glycan alphabet is not universal: a hypothesis**

Jaya Srivastava (ORCID: 0000-0002-1657-4004)[1]*, P. Sunthar[2] and Petety V. Balaji (ORCID: 0000-0002-6018-6957)[1]

[1]Department of Biosciences and Bioengineering, Indian Institute of Technology Bombay, Powai, Mumbai 400076, India

[2]Department of Chemical Engineering, Indian Institute of Technology Bombay, Powai, Mumbai 400076, India

*Corresponding author

Email: jaya_srivastava@iitb.ac.in

Keywords: Glycobiology, bioinformatics, data mining

## Abstract

Several monosaccharides constitute naturally occurring glycans but it is uncertain if they constitute a universal set like the alphabets of proteins and DNA. Based on the available experimental observations, it is hypothesized herein that the glycan alphabet is not universal. Data on the presence / absence of pathways for the biosynthesis of 55 monosaccharides in 12939 completely sequenced archaeal and bacterial genomes are presented in support of this hypothesis. Pathways were identified by searching for homologs of biosynthesis pathway enzymes. Substantial variations are observed in the set of monosaccharides used by organisms belonging to the same phylum, genera and even species. Monosaccharides are grouped as Common, Less Common and Rare based on their prevalence in Archaea and Bacteria. It is observed that fewer enzymes suffice to biosynthesize the Common group. It appears that the Common group originated before the formation of three domains of life. In contrast, the Rare group are confined to a few species in a few phyla, suggesting that they evolved much later. Fold conservation, as observed in aminotransferases and SDR superfamily members involved in monosaccharide biosynthesis, suggests neo- and sub-functionalization of genes leading to the formation of Rare group monosaccharides. Non-universality of the glycan alphabet begets questions about the role of different monosaccharides in determining an organism's fitness.

**Impact statement**

Carbohydrates, nucleic acids and proteins are important classes of biological macromolecules. The universality of DNA, RNA and protein alphabets has been established beyond doubt. However, the universality of glycan alphabet is unknown primarily because of the challenges associated with the elucidation of glycan structures. This has precluded a comprehensive investigation of glycan alphabet. To address this challenge, we have identified the prevalence of 55 monosaccharide biosynthesis pathways in 12939 completely sequenced archaeal and bacterial genomes by searching for homologs of biosynthesis pathway enzymes using HMM profiles, and in a few cases, BLASTp. This revealed that the glycan alphabet is highly variable; in fact, significant differences are found even among different strains of a species. Possible implications of this variability may be significant in understanding the evolution of Archaea and Bacteria in diverse and competitive environments. Factors that drive the choice of monosaccharides used by an organism need to be investigated, and will be of interest in understanding host-pathogen interactions. Additionally, the knowledge of glycan alphabet can be employed for structural characterization / validation of glycans inferred using mass spectrometry. Knowledge of unique monosaccharides and biosynthetic enzymes can also be used as novel drug targets against human pathogens.

**Data summary**

The curated set of proteins used in this study, with domain assignment, is listed in supplementary_data.xlsx. Corresponding 396 references with evidence of experimental characterization are included in supplementary material. Results of genome scan which include predictions of monosaccharides as well as the biosynthesis pathway enzymes is available at http://www.bio.iitb.ac.in/glycopathdb/ including the aforementioned information. Python script used to scan genomes to search for monosaccharide biosynthesis pathways are available on request.

## Introduction

Living organisms show enormous diversity in organization, size, morphology, habitat, etc., but are unified by the highly conserved processes of central dogma: replication, transcription and translation. The enormous diversity seen in life forms is encoded by DNA and decoded primarily by proteins. Both DNA and proteins use the same set of building blocks (nucleotide bases and amino acids, respectively) in all organisms; yet, they store the requisite information by merely varying the (i) set/subset of building blocks used, (ii) number of times each building block is used and (iii) sequence in which the building blocks are linked [collectively referred to as the 'sequence' (Table 1)]. The information required for several other biological processes are stored by glycans, the third group of biological macromolecules (1). It has been found that glycans evolve rapidly in response to changing environmental conditions, especially in Bacteria, and thus contribute to organismal diversity (2,3). The question is, do glycans use the same set of building blocks (viz., monosaccharides) in all organisms, the way proteins and nucleic acids do?

Monosaccharides show a lot more structural variation than amino acids in terms of the enantiomeric forms (both D and L), size (5 to 9 carbon atoms), ring type (pyranose, furanose), and type and extent of modification (deoxy, amino, N-formyl, N-acetyl, etc.). Some pairs of monosaccharides differ from each other merely in the configuration of carbon atoms. The sequence [as defined above] of monosaccharides brings about diversity even in the primary structure of glycans. DNA and protein are linear polymers and the linkage type that connects monomers remains the same throughout. In contrast, glycans can be branched and have alternative isomeric linkages (e.g., $\alpha1\rightarrow3$, $\beta1\rightarrow4$, $\alpha2\rightarrow6$ and so on) (4), two features that enhance diversity in glycans. Repeat length heterogeneity (the number of occurrences of a sequence repeat) is observed in glycans (5,6), as well as DNA and proteins, although there are no data on the frequency of occurrence of this feature in these three classes of biomolecules. An additional factor that contributes to the diversity in the primary structure of glycans is microheterogeneity (7), a feature not seen in DNA or proteins (Table 1). These structural variations demand the use of multiple analytical techniques for sequencing and hence there are no

4

98 automated methods for sequencing glycans. Biosynthesis of DNA and proteins is
99 template-driven but not that of glycans. Consequently, there is no equivalent of
100 polymerase chain reaction or recombinant protein expression to 'amplify' glycans
101 to obtain samples in amounts required for structural / functional analysis. These
102 constraints have largely limited data on glycan sequences.

103

104 **Table 1 Sources of diversity in primary structures of DNA, proteins and glycans**

| Feature | DNA | Protein | Glycan |
|---|---|---|---|
| *Structural diversity of building blocks* | (i) Low (four nucleotides). (ii) Nucleotide modifications (known but rare): N7-methylation of Ade/Gua | (i) Higher, relative to DNA (20 amino acids). Has structurally similar pairs: Asp/Glu, Asn/Gln, Phe/Tyr, Leu/Ile/Val. (ii) Amino acid modifications (known but rare): hydroxylation of Pro and Lys, selenocysteine, pyrrolysine | (i) Highest. Several pentoses and hexoses, many of which are configurational isomers. (ii) Pyranose and furanose forms (e.g., Gal). (iii) Both enantiomeric forms (e.g., Gal). (iv) Modifications extremely common (deoxy, uronic acid, deoxyamino and its derivatives, acetylation, sulfation, …) |
| *Linkage* | 3',5'-phosphodiester. 5',5'-phosphodiester occurs but very rare | Amide bond. $\gamma$-COOH of Glu and $\varepsilon$-NH2 group of Lys used but very rare | Alternative isomeric linkages are very common ($\alpha1\rightarrow3$, $\beta1\rightarrow3$, $\alpha1\rightarrow6$, $\beta1\rightarrow4$, $\alpha2\rightarrow3$, $\alpha2\rightarrow6$ and so on) |
| *Sequence* | (i) Set/subset of building blocks used (ii) Number of times each building block is used (iii) Sequence in which the building blocks are linked | | |
| *Branching* | Absent | Absent | Quite common |
| *Sequence repeat heterogeneity* | Present | Present | Present |
| *Microheterogeneity[1]* | Absent | Absent | Present |

105 [1]Microheterogeneity refers to the presence of multiple forms of glycans (minor but distinct variations)
106 present in different molecules of a protein synthesized by a cell at the 'same' time. This feature is unique
107 to glycans just as the presence of splice variants is unique to proteins.

108

5

109  Monosaccharides are viewed as the third alphabet of life (8). How large is this
110  alphabet? The number of monosaccharides used collectively by living systems is at
111  least 60. An analysis of the bacterial glycan structural data showed a distinct
112  difference in the set of monosaccharides used by bacteria and mammals (9). Is this
113  difference evidence of absence i.e., monosaccharides found in databases are true
114  representations of monosaccharides used by these organisms, and those not found
115  are not used by organisms? Or, is it just absence of evidence i.e., the glycan
116  alphabet is indeed universal and the observed differences are merely due to
117  inadequate sequencing? With the availability of the whole genome sequence of a
118  large number of organisms, it has now become possible to resolve this issue.

119

120  In this study, it is hypothesized that the glycan alphabet is NOT universal, i.e.,
121  different organisms use different sets of monosaccharides. This is in contrast to
122  those of DNA, RNA and proteins. This hypothesis is put forward based on the
123  observations that >60 monosaccharides are found in living systems; the database of
124  glycan structures shows differential usage of monosaccharides and that several
125  serotypes differ from each other in the monosaccharides they use. Results obtained
126  by mining whole genome sequences of 303 Archaea and 12636 Bacteria are
127  presented herein in support of this hypothesis. Monosaccharides considered in this
128  study are nucleotide activated moieties which are utilized by glycosyltransferases
129  (GTs) in the biosynthesis of glycans. Subsequent to such a GT-catalysed transfer,
130  monosaccharides may be modified (e.g., O-acetylation). Monosaccharide
131  derivatives so obtained are not considered in the present study. Enzymes catalysing
132  one or more steps of the biosynthesis pathway are not characterized experimentally
133  for some of the monosaccharides. Such monosaccharides were not considered in
134  this study.

135

136  **Methods**
137  **Databases and software:** Protein sequences and 3D structures were obtained from
138  UniProt and PDB (Table S1). Completely sequenced genomes of 303 Archaea and
139  12636 Bacteria were obtained from the NCBI RefSeq database. These genomes are
140  spread across 3384 species belonging to 1194 genera (Figure S1). Gene
141  neighborhood was analyzed using feature tables taken from NCBI for the
142  respective genomes. BLASTp, MUSCLE, HMMER and CD-Hit (Table S1) were

6

143  installed and used locally. Default values were used for all parameters except when
144  stated otherwise. Word size was set to 2 for BLASTp to prioritize global
145  alignments over local alignments. Thresholds for Hidden Markov Model (HMM)
146  profiles were set based on the best 1 domain bit score rather than e-values since the
147  former is independent of database size.

148  **Searching genomes for monosaccharide biosynthesis pathways**: Pathways for
149  the biosynthesis of 55 monosaccharides have been elucidated to date (Table 2,
150  Figure S2). HMM profiles were generated using carefully curated sets of homologs
151  for 57 families of enzymes that catalyze various steps of 55 monosaccharides
152  (Supplementary_data.xlsx:Worksheet1). Sequences were used directly as BLASTp
153  queries when the number of enzymes characterized experimentally is not sufficient
154  for a HMM profile (Supplementary_data.xlsx:Worksheet2). In-house python
155  scripts were used to scan genomes to identify homologs. Presence of a homolog for
156  each and every enzyme of the biosynthetic pathway of a monosaccharide is taken
157  as evidence of the utilization of this monosaccharide by the organism. On the other
158  hand, absence of a homolog for even one enzyme of the pathway is interpreted as
159  the absence of the corresponding monosaccharide from the organism's glycan
160  alphabet.

161  **Choice of precursors**: Glucose-1-phosphate, fructofuranose-6-phosphate and
162  sedoheptulose-7-phosphate are precursors for many of the monosaccharides
163  (Supplementary_data.xlsx:Worksheet6).      Fructofuranose-6-phosphate      and
164  sedoheptulose-7-phosphate are intermediates in the glycolytic pathways viz.,
165  Embden-Meyerhof pathway and pentose phosphate pathway, respectively, and
166  these enzymes are not considered for the search. Pathways for biosynthesis of
167  UDP-Glc2NAc and GDP-mannose have been considered separately since
168  Glc2NAc and mannose are glycan building blocks as well as intermediates in the
169  biosynthesis of several other monosaccharides. Hence biosynthesis steps of UDP-
170  Glc2NAc and GDP-mannose were excluded from those of their derivatives. An
171  additional pathway for UDP-glucose biosynthesis was considered to analyze its
172  ubiquity since UDP-glucose is part of both anabolic and catabolic pathways. The
173  biosynthesis of CMP-Leg5Ac7Ac starting from N-acetyl-glucosamine-1-phosphate

174 has also been considered because of the uncommon guanylyltransferase in the first
175 step of the pathway.

176 **Table 2 Summary of the pathways for the biosynthesis of monosaccharides[a,b]**

| Details about the end product of biosynthesis pathways | Precursor[c] | | | | | |
|---|---|---|---|---|---|---|
| | Glc-1-P | Fru$f$-6-P | GDP-Man | UDP-Glc2NAc | Glc2NAc-1-P | Sed-7-P |
| Number of nucleotide sugars[d] | 27 | 2 | 8 | 16 | 1 | 4 |
| Number of monosaccharides[e] | 25 | 2 | 8 | 16 | 1 | 4 |
| Number of monosaccharides with different number of backbone carbon atoms | | | | | | |
| Pentose | 4 | - | - | - | - | - |
| Hexose | 21 | 2 | 8 | 13 | - | - |
| Heptulose | - | - | - | - | - | 4 |
| Nonulose | - | - | - | 3 | 1 | - |
| Number of monosaccharides of the two enantiomeric forms[f] | | | | | | |
| D | 19 | 2 | 5 | 12 | 1 | 3 |
| L | 6 | - | 3 | 4 | - | 1 |
| Number of monosaccharides of the two ring forms | | | | | | |
| Pyranose | 23 | 2 | 8 | 16 | 1 | 4 |
| Furanose | 2 | - | - | - | - | - |
| Number of monosaccharides with different nucleotides | | | | | | |
| ADP | - | - | - | - | - | 1 |
| CDP | 7 | - | - | - | - | - |
| CMP | - | - | - | 3 | 1 | - |
| GDP | - | 1 | 8 | - | - | 3 |
| TDP/dTDP[g] | 9 | - | - | - | - | - |
| UDP | 11 | 1 | - | 13 | - | - |

177
178 [a] The monosaccharide L-Iduronic acid has not been considered in this study since there is no separate pathway for its
179 biosynthesis. Dermatan sulfate epimerase-1 or -2 (DS-epi1 or DS-epi2) catalyses C5-epimerization of glucuronic
180 acid to L-iduronic acid in chondroitin sulfate polymeric chains (10).
181 [b] Enzymes catalysing one or more steps of the biosynthesis pathway are not characterized experimentally for some
182 of the monosaccharides. Such monosaccharides were not considered in this study.
183 [c] Glc-6-P is the precursor for Glc-1-P (conversion catalysed by phosphoglucomutase), Fru$f$-6-P (catalysed by
184 phosphoglucose isomerase) and Sed-7-P (formed in the non-oxidative phase of the pentose phosphate pathway).
185 Fruf-6-P is the precursor of GDP-Man and UDP-Glc2NAc.
186 [d] There are two pathways for the biosynthesis of CMP-Leg5Ac7Ac, one starting from UDP-Glc2NAc and the other
187 from Glc2NAc-1-P. Hence, the total number of nucleotide sugars will be 57 even though row sum is 58.
188 [e] L-Rhamnose and Qui4NAc are biosynthesized as both UDP- and TDP-/dTDP-derivatives. Hence, the number of
189 monosaccharides is less than the number of nucleotide sugars by 2.
190 [f] The prefix D is omitted for D enantiomers whereas the prefix L is explicitly mentioned for L enantiomers.

191   [g] No distinction is made between TDP and dTDP in this work since literature suggests that both ribo- and deoxyribo-
192   substrates are used by enzymes, albeit with varying extents of specificity depending upon the source organism. In
193   fact, dTDP and TDP have been used synonymously by some authors.

194

195   **Generation of HMM profiles**: An HMM profile was generated for each step of a
196   biosynthesis pathway except where mentioned otherwise. Profiles were generated
197   in two steps (Flowchart S1). The extended dataset was created to account for
198   sequence divergence. In some cases, no additional sequences satisfying the
199   aforementioned criteria were found, hence there is no Extend dataset. Each profile
200   was given an annotation based on the enzyme activities of proteins that were used
201   to generate the profile and an identifier of the format GPExxxxx; here GPE stands
202   for Glycosylation Pathway Enzyme and xxxxx is a unique 5-digit number
203   (Supplementary_data.xlsx:Worksheet1).

204

205   **Setting thresholds for HMM profiles**: Thresholds for HMM profiles were set as
206   described        below        (profile-wise        details        are        given        in
207   Supplementary_data.xlsx:Worksheet1):

208

209   Using ROC curves: TrEMBL database was used to generate ROC curves. Several
210   of the TrEMBL entries have been assigned molecular function electronically based
211   on UniRule and SAAS (Table S1). It is assumed that these annotations are correct
212   while generating ROC curves. True positives, false positives and false negatives
213   were identified by comparing TrEMBL annotation with profile annotation.

214

215   Using bit-score scatter plots: Members of some enzyme families differ in their
216   molecular function while retaining significant global sequence similarity e.g., C4-
217   and C3-aminotransferases. Consequently, annotations of several TrEMBL
218   sequences belonging to such families are incomplete e.g., DegT/DnrJ/EryC1/StrS
219   aminotransferase family protein. In such cases, bit score scatter-plots were used to
220   set thresholds (Figure S5). Scatter plot was also used to set threshold in case of
221   hydrolysing and non-hydrolysing NDP-Hex2NAc C2 epimerases since many
222   TrEMBL hits are just annotated as NDP-Hex2NAc C2 epimerases.

223

224　Using $T_{exp}$ and $T_{extend}$ as thresholds: $T_{exp}$ or $T_{extend}$ was used as the threshold for
225　some profiles for one of these two reasons: (i) Sequences used to generate the
226　profile are a subset of the sequences used to generate another profile; the latter set
227　of enzymes has broader substrate specificity than those of the former set. For
228　instance, sequences used for generating GPE02430 [TDP-/dTDP-4-keto-6-
229　deoxyglucose 3-/3,5-epimerase] and GPE02530 (NDP-sugar 3-/3,5-/5-epimerase)
230　are homologs but the former set has narrow specificity. $T_{extend}$ was set as
231　threshold for GPE02430 as lowering the threshold would make this profile less
232　specific. (ii) For some profiles such as GPE50010 [nucleotide sugar
233　formyltransferase], very few TrEMBL entries that score $< T_{exp}$ have been assigned
234　molecular function and hence ROC curve could not be generated.

235

236　The case of GPE00530: Scanning TrEMBL database with GPE00530 (Glucose-1-
237　phosphate uridylyltransferase family 2) using the default threshold of HMMER (e-
238　value = 10) resulted in 2693 hits with matching annotation and their scores ranged
239　continuously from 705 to 303 bits and then from 57 to 41 bits. It was not possible
240　to generate a ROC curve because of this discontinuity. Hence, 303 bits was set as
241　the threshold.

242

243　**Profile annotations with broader substrate / product specificities**: Many
244　sequence homologs catalyse the "same" reaction but with (slightly) different
245　substrate specificities. Sequence changes that confer such differential specificities
246　are subtle and often unknown. HMM profiles of such families lack the ability to
247　discriminate between sequences with varying substrate specificities. Two products,
248　a major product and a minor product, are formed in certain enzyme catalysed
249　reactions (11–13). It is possible that only the major product has been characterized
250　while assaying an enzyme with broader substrate specificity. Another possibility is
251　that only a subset of possible substrates has been assayed for. Hence, substrate
252　specificities are broad in annotations of some of the profiles. As opposed to these,
253　some profiles of aminotransferases and reductases are generated from enzymes
254　which differ from each other with respect to the product formed viz., orientation
255　(equatorial or axial) of the newly formed/added -OH / -NH2 group. Profile for 3,4-
256　ketoisomerase is also of this type. UDP-GlcA decarboxylase (UXS) converts UDP-

10

257 GlcA to UDP-4-keto xylose, which is further reduced to UDP-xylose. UDP-4-keto
258 xylose is a minor product for human UXS whereas it is a major product for *E. coli*
259 UXS (12). Both these enzymes are used to generate the profile GPE20030
260 (Supplementary_data.xlsx:Worksheet1).

261

262 <u>Pathway steps associated with more than one HMM profile</u>: Some steps are
263 associated with more than one profile for one of these two reasons: (i) Non-
264 orthologous enzymes known to catalyse the same reaction e.g.,
265 phosphomannoisomerases. (ii) Two or more profiles are generated, one with
266 narrow and the other(s) with broad substrate specificity. Enzymes used for the
267 former are a subset of enzymes used for the latter type of profiles e.g.,
268 aminotransferases. The process flow adopted to assign annotation for a sequence
269 which satisfies thresholds for more than one profile is shown in Flowchart S2.

270

271 **Finding homologs using BLASTp instead of HMM profiles**: HMM profiles
272 were generated only when four or more experimentally characterized enzymes are
273 available (two exceptions are discussed below). Global alignment and sequence
274 similarity were used as the criteria to infer homology based on BLASTp search.
275 The default values were set to be $>= 90\%$ query coverage and $>=30\%$ sequence
276 similarity. However, these values were upwardly revised when query sequences
277 belonged to homologous families that are functionally divergent
278 (Supplementary_data.xlsx:Worksheet2). Specifically, similarity and coverage cut-
279 offs were revised by performing an all-against-all BLASTp search of all
280 experimentally characterized sequences of monosaccharide biosynthesis pathways.

281

282 *B. cereus* PdeG (Q81A42_1-328) is a retaining UDP-Glc2NAc 4,6-dehydratase
283 (14). It shares higher sequence similarity with inverting UDP-Glc2NAc 4,6-
284 dehydratases than with retaining dehydratases. The sequence of PdeG was
285 compared with TrEMBL hits for the HMM profile of inverting UDP-Glc2NAc 4,6-
286 dehydratases (GPE05331), based on which the sequence similarity cut-off for
287 PdeG was set to 70%. The threshold for GPE05331 was set such as to exclude
288 PdeG (Figure S5).

289 Criteria for finding homologs of UDP-2,4-diacetamido-2,4,6-trideoxy-β-L-altrose
290 hydrolase and UDP-4-amino-6-deoxy-Glc2NAc acetyltransferase: Four
291 experimentally characterized enzymes are known for each of these two families.
292 However, BLASTp approach was used instead of generating an HMM profile. This
293 is because a suitable bit score threshold could not be arrived, which in turn, was
294 because several of the TrEMBL entries obtained as hits are annotated as CMP-N-
295 acetylneuraminic acid synthetase or equivalent (for hydrolase), or O-
296 acetyltransferase or equivalent (for acetyltransferase).

297 **Uncertainties in prediction:** Any description of molecular function of a protein is
298 stratified and includes specifying the type of reaction catalysed, substrate(s) used,
299 etc. A vast majority of sequences conceptually translated from genome sequences
300 are assigned molecular function based on sequence homology to experimentally
301 characterized proteins. Even though experimental validation is available for only a
302 small fraction of proteins due to practical constraints, such studies have shown that
303 homology-based assignments are generally valid and deviations typically pertain to
304 the extent of substrate specificity, metal ion dependency and such. Nevertheless,
305 caution is warranted with increasing sequence divergence and one has to be on the
306 lookout for homologs that have acquired new molecular function as a result of
307 mutation of a handful of key residues (neo-functionalization). In view of this, in
308 the present study, HMM and BLASTp thresholds have been chosen with higher
309 stringency and assignment of substrate(s) and product(s) has been made
310 conservatively by manually curating false positives and false negatives from the
311 Swiss-Prot database, details of which are given below:

312

313 (1) Both GDP-rhamnose and GDP-6-deoxytalose are assigned as products of the
314 same pathway, because their biosynthesis proceeds through the same pathway with
315 the exception of the last step being catalysed by homologous 4-reductases. It is not
316 possible to infer if product specificity of enzymes in this family is absolute or
317 partial i.e., one is a major product and other, a minor product, due to inadequate
318 experimental data. An identical situation is seen in the pathways for the
319 biosynthesis of CDP-cillose and CDP-cereose, and for CDP-abequose and CDP-

paratose. In view of this, prevalence data will be the same for the two monosaccharides of a pair (Supplementary_data.xlsx:Worksheet3).

(2) Non-hydrolyzing NDP-Hex2NAc C2-epimerases (GPE02030) are part of biosynthesis pathways of different monosaccharides. The extent of substrate specificity of the experimentally characterized members of this family is not known since not all enzymes have been assayed using all possible substrates. In literature, substrate specificity is arrived at based on the genomic context and the same approach has been followed in the present study as well. For example, hits for GPE02030 profile are treated as Man2NAc synthesis pathway enzymes, unless other enzymes of L-Fuc2NAc, L-Qui2NAc or Man2NAc3NAcA pathway are also present.

(3) Some monosaccharides are precursors for other monosaccharides and hence, genomes predicted to have the pathway for the latter monosaccharide will also have the precursor monosaccharide. Following are the precursor-final product monosaccharide pairs encountered in this study: (i) L-Rha2NAc $\rightarrow$ L-Qui2NAc, (ii) L-Rhamnose $\rightarrow$ 6-Deoxy-L-talose, (iii) Fucose $\rightarrow$ Fucofuranose, (iv) Paratose $\rightarrow$ Tyvelose, (v) Galactose $\rightarrow$ Galactofuranose, (vi) GlcA $\rightarrow$ GalA, (vii) L-Ara4N $\rightarrow$ L-Ara4NFo, (viii) Per $\rightarrow$ Per4Ac, (ix) Man2NAc $\rightarrow$ Man2NAcA, (x) Glc2NAcA $\rightarrow$ Gal2NAcA, and (xi) Bac2Ac4Ac $\rightarrow$ Leg5Ac7Ac.

(4) The pathway for the synthesis of L-arabinose is an extension of the pathway for the synthesis of xylose. However, most genomes predicted to have xylose pathway also have L-arabinose pathway. This is because UDP-sugar C4-epimerase family members (GPE02230) catalyse C4-epimerization of glucose, GlcA, Glc2NAc, Glc2NacA and xylose. Assigning substrate specificity solely based on sequence similarity is not possible. The challenge is compounded by the fact that some of these enzymes show broad substrate specificity while the rest are only specific to a single substrate. Not all enzymes have been assayed for all potential substrates.

**Results**

**Glycan alphabet size is not the same across Archaea and across Bacteria**: The number of monosaccharides used by different species is significantly different (Figure 1) and is independent of proteome size (Figure S3). Data for the prevalence of monosaccharides in 12939 genomes is very similar to that in 3384 species (Figure S4) indicating that the outcome is not biased by the skew in the number of genomes (strains) sequenced for a given species (Figure S1). In fact, none of the organisms use all 55 monosaccharides: the highest number of monosaccharides used by an organism is 23 [*Escherichia coli* 14EC033]. Just 1 and 2 monosaccharides are used by 188 and 117 species, respectively. Glucose, galactose and mannose, and their 2-N-acetyl (Glc2NAc, Gal2NAc, Man2NAc) and uronic acid (GlcA, GalA, Glc2NAcA, Gal2NAcA) derivatives are the most prevalent besides L-rhamnose, as the biosynthesis pathways for these monosaccharides are found in >50% of genomes (Figure 2). These monosaccharides are thus categorized as 'Common' group. However, none of them are used by all organisms (Supplementary_data.xlsx:Worksheet3).

**Figure 1 The number of monosaccharides for which biosynthesis pathways are found in a species**. More than one strain is sequenced for several species (Figure S1). In such cases, data for the strain which has the highest number of monosaccharides is plotted. Total number of species = 3384.

**Figure 2 Classification of monosaccharides into three groups based on prevalence in archaeal+bacterial genomes**. These groups are Common (found in >=50% of genomes), Less Common and Rare (found in <=10% of genomes). Abbreviated names are used for some of the monosaccharides. Full names of these are given in Supplementary_data.xlsx:Worksheet4.

**Evolution and diversification of glycan alphabet**: It is observed that only a limited set of enzymes suffice to biosynthesize the Common group monosaccharides e.g., nucleotidyltransferases (activation), amidotransferase and N-acetyltransferase (Hex2NAc from a hexose), C4-epimerase (Glc to Gal) and C6-dehydrogenase (uronic acid) belonging to the SDR superfamily, non-hydrolyzing C2-epimerase (Glc2NAc to Man2NAc), mutase (6-P to 1-P) and isomerase (pyranose to furanose) (Figure 3 and Supplementary_data.xlsx:Worksheet5). Using this limited set of monosaccharides, organisms seem to achieve structural diversity

14

382 by mechanisms such as alternative isomeric linkages, branching and repeat length
383 heterogeneity. Some organisms use an additional set of monosaccharides, viz., L-
384 fucose, galactofuranose, xylose, L-Ara4N and L-arabinose. These monosaccharides
385 are categorized as Less Common group. Organisms using this group of
386 monosaccharides have enhanced the glycan repertoire by acquiring C3/C5-
387 epimerase, 4,6-dehydratase, C4-reductase, C6-decarboxylase and C4-
388 aminotransferase. The rest of the monosaccharides are used by very few organisms
389 and thus constitute the 'Rare' group (Figure 2).

390

391 **Figure 3 A qualitative comparison of the number of monosaccharides of the three groups viz.,**
392 **Common (C), Less Common (LC) and Rare (R) with their prevalence in archaeal+bacterial**
393 **genomes and the number of types of enzymes required for their biosynthesis**. The size of a group is
394 inversely related to the prevalence of the corresponding group of monosaccharides. Enzymes required for
395 the biosynthesis of Common group monosaccharides are required for the biosynthesis of Less Common
396 and Rare groups also; similarity, those for the Less Common group are required for the biosynthesis of
397 Rare group also. Different enzymes belonging to each of the superfamily mentioned above are listed in
398 the file Supplementary_data.xlsx:Worksheet5. Note that the group sizes are not to scale. It may be noted
399 that additional types of enzymes may have to be included when experimental data about the pathways for
400 the biosynthesis of other monosaccharides becomes available.
401 HAD, haloalkanoic acid dehalogenase
402 Gfo/Idh/MocA, glucose☐fructose oxidoreductase/inositol 2☐dehydrogenase/rhizopine catabolism protein
403 MocA
404 GNAT, GCN5-related N-acetyltransferases
405 LβH, left handed β helix
406 PEP, phosphoenolpyruvate
407 PLP, pyridoxal 5'-phosphate
408 SAM, S-adenosyl-L-methionine
409 SDR, short chain dehydrogenase reductase
410 UDP, uridine diphosphate

411

412 Occurrence of the Common group of monosaccharides in all three domains of life
413 points to their presence early on during evolution. Neo- and sub-functionalization
414 of horizontally acquired and duplicated genes during the course of evolution have
415 been widely reported (e.g.,(15,16)). It is envisaged that the enzymes required for
416 the biosynthesis of Rare group monosaccharides have arisen by such neo- and sub-
417 functionalization. Aminotransferase and short-chain dehydrogenase reductase
418 (SDR) superfamily enzymes involved in the biosynthesis of monosaccharides lend
419 support to this inference. Superposition of a few C3- and C4-aminotransferases

15

420  show remarkable conservation of the 3D structures despite differences in the
421  pyranose ring position at which the amino group is transferred as well as the
422  nucleotide sugar substrate (Figure 4). 3D structures are conserved even among
423  SDR superfamily enzymes despite catalysing different reactions viz., epimerization
424  (at C2 or C4), removal of water (dehydratase at C4, C6) and reduction (at C4).

425

426  **Figure 4 3D structural superimposition of enzymes belonging to aminotransferase (A) and SDR (B)**
427  **superfamilies involved in the biosynthesis of monosaccharides**. Color scheme: helices, raspberry red;
428  sheet, forest green; loops, light blue. **Panel A**: Aminotransferase superfamily enzymes: 1MDO_A: ArnB
429  from UDP-L-Ara4N biosynthesis; 2FNI_A: PseC from CMP-L-Pse45Ac7Ac biosynthesis; 2OGA_A:
430  DesV from TDP-/dTDP-desosamine biosynthesis; 3BN1_A: perA from GDP-per biosynthesis; 3NYU_A:
431  WbpE from UDP-Man2NAc3NAcA biosynthesis; 4PIW_A: WecE from TDP-/dTDP-Fuc4NAc
432  biosynthesis; 4ZTC_A: PglE from CMP-Leg5Ac7Ac biosynthesis; 5U1Z_A: wlaRG from TDP-/dTDP-
433  Fuc3NAc/Qui3NAc biosynthesis. ArnB, PseC, perA, WecE and PglE are C4-aminotransferases whereas
434  DesV, WbpE and wlaRG are C3-aminotransferases. **Panel B**: 1ORR_A: RfbE, C2-epimerase from CDP-
435  tyvelose biosynthesis; 2PK3_A: Rmd, 4-reductase from GDP-rhamnose biosynthesis; 1KBZ_A: rmlD,
436  C4-reductase from TDP-/dTDP-L-rhamnose biosynthesis; 1T2A_A: gmd, C4,C6-dehydratase from GDP-
437  L-fucose biosynthesis; 1SB8_A: WbpP, C4-epimerase from UDP-Gal2NAc biosynthesis; 5BJU_A: PglF,
438  C4,C6-dehydratase from UDP-Bac2Ac4Ac biosynthesis.

439

440

441  **Glycan alphabet varies even across strains**: Remarkably, variations in the size of
442  glycan alphabet are significant even at the strain level (Figure 5). Strain-specific
443  differences are pronounced in species such as *E. coli, Pseudomonas aeruginosa*
444  and *Campylobacter jejuni* (Figure 6) possibly reflecting the diverse environments
445  that these organisms inhabit. Among organisms which inhabit the same
446  environment, strain-specific differences show mixed pattern: among the 71 strains
447  of *Streptococcus pneumoniae*, the maximum and minimum number of
448  monosaccharides utilized by a strain are 4 and 12, respectively. Such a variation
449  could have evolved as a mechanism to evade host immune response. In contrast,
450  strains of *Streptococcus pyogenes* and strains of *Staphylococcus aureus* inhabit the
451  same environment (respiratory tract and skin, respectively) and show very little
452  variation in the monosaccharides they use. Both are capsule producing
453  opportunistic pathogens suggesting that they might bring about antigenic variation
454  by variations in linkage types, branching etc. (17), even with the same set of
455  monosaccharides. Strains of *Mycobacterium tuberculosis, Brucella melitensis,*
456  *Brucella abortus* or *Neisseria gonorrhoeae*, all of which are human intracellular
457  pathogens, also show insignificant variation. It is possible that different strains of a

458 pathogen are a part of distinct microbiomes and microbial interactions within the
459 biome/with the host determine the glycan alphabet of the organism. Availability of
460 additional characteristics such as phenotypic data and temporal variations in glycan
461 structures is critical for understanding the presence/absence of strain-specific
462 variations.
463

464 **Figure 5 Variations in the number of monosaccharides used by different strains of a species.** Species
465 with more than one sequenced strain and at least one monosaccharide predicted in one of the strains are
466 considered. Only the smallest and largest numbers are shown.

467 **Figure 6 Different strains of some of the species do not use the same number of monosaccharides**.
468 The range of the number of monosaccharides used by various strains of some of the clinically important
469 species are shown here. The number of sequenced strains for each organism is shown above the
470 corresponding bar. Number in parenthesis after the name of each organism represents the minimum
471 number of monosaccharides used by one of the strains of this organism. Note that the set of
472 monosaccharides encoded by different strains utilizing the same number of monosaccharides may vary.
473 Organisms associated with narrow habitat are shown in blue, while those with broad habitat are shown in
474 purple.

475 **Prevalence of monosaccharides across Phyla**: Not all sugars of the Common
476 group (Figure 1) are found across all phyla whereas Neu5Ac belonging to the Rare
477 group is found across all phyla. GlcA and GalA (Common group) are absent in
478 Thermotogae suggesting that pathways for their biosynthesis are lost in this
479 phylum. A similar conclusion is drawn for the absence of L-fucose and L-colitose
480 in TACK group phylum. Most of the Rare group sugars are limited to very few
481 species in a few phyla (Figure 7). For instance, Fuc4NAc and L-*glycero*-β-D-
482 *manno*-heptose (ADP-linked) are found only in Gamma-proteobacteria, a class that
483 comprises of several pathogens. The other three heptoses, which are GDP-linked,
484 are absent in Gamma-proteobacteria. Recently, it was found that *Helicobacter*
485 *pylori*, belonging to the class Epsilon-proteobacteria, synthesizes ADP-*glycero*-β-
486 D-*manno*-heptose for activating the NF-κβ pathway in human epithelial cells (18).
487 This pathway has been experimentally characterized in very few organisms.
488 Consequently, homologs for this pathway are found by BLASTp queries and not
489 by HMM profiles. In the present study, this pathway turned out to be a false
490 negative because of the high stringency set for BLASTp thresholds. In view of this,
491 it is possible that such sugars which appear restricted to a few phyla are also found
492 in others.
493

494 **Figure 7 Prevalence of Less Common and Rare group monosaccharides in different microbial**
495 **phyla.** Data for phyla with less than five sequenced genomes are not shown to avoid visual clutter. Only
496 names of monosaccharides are used for annotation even though all are biosynthesized as nucleotide
497 sugars. Abbreviated names are used for some of the monosaccharides. Full names of these are given in
498 Supplementary_data.xlsx:Worksheet4.
499

500 **Why some eubacteria do not biosynthesize any monosaccharide?** None of the
501 monosaccharides are biosynthesized by some mollicutes (e.g., Mycoplasma) and
502 endosymbionts (e.g., *Ehlrichia sp.* and *Orientia sp.*) because the biosynthesis
503 pathways are completely absent. Mollicutes lack cell wall (19) which could explain
504 the absence of monosaccharides. Endosymbionts have reduced genomes which is
505 seen as an adaptation to host dependence (20) (21). Biosynthesis pathway enzymes
506 are lost / are being lost as part of the phenomenon of genome reduction. This is
507 illustrated by the endosymbiont *Buchnera aphidicola*: 13 of the 25 strains have the
508 pathway for the biosynthesis of UDP-Glc2NAc, 7 have partial pathway and 5 do
509 not encode any gene of this pathway. Pathway for none of the other
510 monosaccharides are found in this organism. Pathways are incomplete i.e.,
511 enzymes catalysing one or more steps of the pathway are absent in some
512 organisms. Some species of Mycoplasma, Ureaplasma and Spiroplasma lack
513 mannose-1-phosphate guanylyltransferase because of which GDP-mannose is not
514 biosynthesized. GlmU which converts Glc2N-1-phosphate to UDP-Glc2NAc is
515 absent in *Chlamydia sp.* However, Glc2N is found in the LPS of *Chlamydia*
516 *trachomatis* (22). Whether this is indicative of the presence of a transferase which
517 uses Glc2N-1-phosphate instead of UDP-Glc2N needs to be explored.

518

519 **Do *Rickettsia sp.* and *Chlamydia sp.* source monosaccharides from their host?**
520 *Rickettsia sp.* (60 strains), *Orientia tsutsugamushi* (7 strains), and *Chlamydia sp.*
521 (143 strains) are obligate intracellular bacteria. *O. tsutsugamushi* does not contain
522 pathways for the biosynthesis of any of the monosaccharides. This is in
523 consonance with the finding that it does not contain extracellular polysaccharides
524 (21). Rickettsia species have pathways for the biosynthesis of Man2NAc, L-
525 Qui2NAc and L-Rha2NAc. L-Rha2NAc is the immediate precursor for L-Qui2NAc
526 (Figure S2e). Rickettsia are known to use Man2NAc and L-Qui2NAc but not L-
527 Rha2NAc (23) implying that UDP-L-Rha2NAc is just an intermediate in these
528 organisms. The pathway for the biosynthesis of UDP-Glc2NAc, precursor for these

529 Hex2NAcs, is absent suggesting partial dependence on host (human). Notably,
530 genes for the biosynthesis of Man2NAc and L-Qui2NAc have so far not been
531 reported in humans, which explains why Rickettsia have retained these pathways
532 (the human genome was scanned and these pathways are not found; unpublished
533 data). Both Rickettsia and Orientia belong to the same order, Rickettsiales.
534 Symptoms caused by these two are similar (24). In spite of similarities in host
535 preference and pathogenicity, *Rickettsia sp.* continues to use certain
536 monosaccharides while diverging from *O. tsutsugamushi* (25) which uses none. Is
537 this because Rickettsia use ticks as vectors whereas Orientia use mites (26)?
538 *Rickettsia akari*, the only Rickettsial species which uses mites as vectors and
539 contains pathways for Man2NAc and L-Qui2NAc biosynthesis, has been proposed
540 to be placed as a separate group because its genotypic and phenotypic
541 characteristics are intermediate to those of Orientia and Rickettsia (26).

542

543 **Absence of Glc2NAc in organisms other than endosymbionts**: UDP-Glc2NAc
544 is the precursor for the biosynthesis of several monosaccharides (Figures S2e, S2f).
545 However, pathways for its biosynthesis are absent in ~10% of the genomes
546 excluding endosymbionts. None of the organisms in FCB group and Spirochaetes
547 contain this monosaccharide. Further analysis revealed the loss of first (GlmS) or
548 last (GlmU) enzyme of the pathway in several of their genomes. This pattern
549 suggests that organisms of this phyla are in the process of losing UDP-Glc2NAc
550 pathway. Incidentally, some of these genomes do contain its derivatives. They
551 include host-associated organisms such as *Bacteriodes fragilis, Flavobacterium*
552 *sp., Tannerella forsythia, Akkermansia muciniphila, Bifidobacterium bifidum,*
553 *Leptospira interrogans*, etc., suggesting that they obtain Glc2NAc from their
554 microenvironment. However, a few free-living organisms which contain
555 derivatives of UDP-Glc2NAc but not UDP-Glc2NAc were also identified. For
556 instance, GlmU is not present in *Arcticibacterium luteifluviistationis* (arctic surface
557 seawater) and its C-terminus (acetyltransferase domain) is absent in
558 *Chlorobaculum limnaeum* (freshwater). Nonetheless, both organisms contain the
559 UDP-L-Qui2NAc pathway cluster.

560

561 **Prevalence of enantiomeric pairs and isomers of N-acetyl derivatives**: Both
562 enantiomers of a few monosaccharides are reported in natural glycans. The two
563 enantiomers may or may not be biosynthesized from the same precursor, and may
564 be linked to different nucleotides (Table S2). The present analysis shows that both
565 enantiomers are found in only a small number of organisms, that too in specific
566 genera, class or phyla (Table 3). Three isomeric N-acetyl derivatives of fucosamine
567 (6-deoxygalactosamine) and of quinovosamine (6-deoxyglucosamine) are found in
568 living systems. The N-acetyl group is present at C2, C3 or C4 position in these
569 isomers. Only few organisms use more than one of these three isomers (Table 3).
570 One such organism is *E. coli* NCTC11151 which contains both Fuc4NAc and
571 Fuc3NAc. In contrast, *E. coli* O177:H21 uses L-Fuc2NAc along with Fuc3NAc.
572 Genomic context analysis showed that Fuc4NAc biosynthesis genes are part of the
573 O-antigen cluster in both these strains. On the other hand, genes for the
574 biosynthesis of Fuc3NAc (in NCTC11151) and L-Fuc2NAc (in O177:H21) are
575 present as part of the colanic acid cluster. Four genomes (strains) of *Pseudomonas*
576 *orientalis* use Qui4NAc, Qui2NAc and L-Qui2NAc; genes required for the
577 biosynthesis of these three monosaccharides are all in the same genomic
578 neighbourhood.

579

580 **Table 3 Presence of enantiomeric pairs and isomeric N-acetyl derivative pairs§**

| Monosaccharide | Where present | Number of organisms (genomes) in which these monosaccharide pairs are used |
|---|---|---|
| **Both D- and L-enantiomers** | | |
| Galactose | Extremophiles | 33 |
| Fucose | Phylum FCB group<br>Phylum Deferribacteres<br>Class Gammaproteobacteria | 7<br>1<br>10 |
| Rhamnose | Genus Pseudomonas | 338 |
| 6-Deoxytalose | Genus Pseudomonas | 140 |
| Qui2NAc[†] | Phylum Proteobacteria | 64 |
| Fuc2NAc | Genus Staphylococcus[¶] | 300 |

| Isomers of N-acetyl derivatives | | |
|---|---|---|
| Fuc3NAc and Fuc4NAc[‡] | Family Enterobacteriaceae | 33 |
| Qui2NAc, Qui4NAc | Several phyla | 193 |

581  [§] The diastereomeric pair of Fuc4NAc and L-Fuc2NAc are found in some strains of *Escherichia coli*.

582  [¶] Both Fuc2NAc and L-Fuc2NAc are components of capsular polysaccharides (27)

583  [†] Abbreviated names are used for some of the monosaccharides. Full names of these are given in
584  Supplementary_data.xlsx:Worksheet4.
585  [‡] Glucose-1-phospate is the precursor for both Fuc3NAc and Fuc4NAc and UDP-Glc2NAc is the
586  precursor for Fuc2NAc (the isomer that is absent in these organisms).

587

## Why are some pathways not found in Archaea?

589  Most of the rare group monosaccharides are absent in Archaea. Members of
590  Euryarchaeota contain higher number of monosaccharides than TACK group. This
591  could be suggestive of lateral gene transfer events with bacterial members as
592  members of Euryarchaeota, particularly methanogens, coexist with other
593  organisms in microbiomes (28) and have been inferred to acquire their genetic
594  content (29). It is premature to associate absence of monosaccharide diversity to
595  the apparent lack of pathogenicity in Archaea (28). This is because of inadequate
596  information regarding the abundance of Archaea in various microbiomes. This in
597  turn is due to our limitations in the detection of Archaea and associating them with
598  disease phenotypes.

599

600  Apart from these possibilities, methodological limitations may have resulted in
601  apparent absence of monosaccharides in Archaea. Only 4-5% of the 789 sequences
602  used for generating HMM profiles or as BLASTp queries are from archaea. The
603  pathway for the biosynthesis of TDP-/dTDP-L-rhamnose has four enzymes viz.,
604  RmlA, RmlB, RmlC and RmlD. Of these, only RmlB could not be found by HMM
605  profile in *Saccharolobus sp., Desulfurococcus sp.* and *Sulfolobus sp.* leading to the
606  conclusion that L-rhamnose is absent in these organisms. Analysis of the
607  neighbourhood of RmlA, RmlC and RmlD revealed a sequence which could
608  potentially be RmlB since it retains conserved residues of this family. This
609  sequence could not be captured by the profile-based search due to stringent
610  thresholds (=400 bits) [profile GPE05430; Supplementary_data.xlsx:Worksheet1].

611 Potential RmlB sequences of these organisms score 300-350 bits. This observation
612 suggests that the pathway exists in these organisms but was not identified due to
613 the stringency of the threshold. However, this is in contrast to other cases of
614 absence of monosaccharides wherein none of the proteins of a pathway in the
615 genome score even the default bit score of HMMER (i.e., 10 bits).

616

617 **Use of more than one nucleotide derivative/alternative pathways:** L-rhamnose
618 and Qui4NAc are biosynthesized as both UDP- and TDP-/dTDP-derivatives
619 (Figures S2a, S2c). However, the TDP-/dTDP-pathways are found in Archaea and
620 Bacteria, but not the UDP-pathways. TDP-/dTDP-6-deoxy-L-talose is
621 biosynthesized via reduction of TDP-/dTDP-4-keto-L-rhamnose or C4
622 epimerization of TDP-/dTDP-L-rhamnose (Figure S2a). The former pathway
623 occurs in 141 genomes belonging to multiple phyla and notably in *Pseudomonas*
624 *sp., Streptococcus sp.* and *Streptomyces sp*. The latter pathway is found in 255
625 genomes belonging to Proteobacteria and Terrabacteria, and notably in
626 *Burkholderia sp., Mycobacterium sp.* and *Xanthomonas oryzae*. N,N'-diacetyl
627 legionaminic acid can be biosynthesized either from UDP-route or GDP-route
628 (Figure S2f). The latter pathway is found in 93 of 96 genomes of *Campylobacter*
629 *jejuni* whereas the former is found in 10 other genomes primarily belonging to
630 Bacteriodetes/Chlorobi class.

631

632 **Discussion**
633 The importance of glycans, especially in Archaea and Bacteria, is well
634 documented. Establishing the specific role of glycans and studying structure-
635 function relationship is largely hindered by factors such as non-availability of high-
636 throughput sequencing methods, inadequate information as to which genes are
637 involved in non-template driven biosynthesis, phase variation (30) and
638 microheterogeneity (7). In this study, completely sequenced archaeal and bacterial
639 genomes were searched for monosaccharide biosynthesis pathways using sequence
640 homology-based approach. It is found that the usage of monosaccharides is not at
641 all conserved across Archaea and Bacteria. This is in stark contrast to the alphabets
642 of DNA and proteins which are universal. In addition, marked differences are
643 observed even among different strains of a species. The range of monosaccharides
644 used by an organism seems to be influenced by environmental factors such as

645 growth (nutrients, pH, temperature, …) and environmental (host, microbiome, …)
646 conditions. For instance, high uronic acid content in exopolysaccharides of marine
647 bacteria imparts anionic property which is implicated in uptake of $Fe^{3+}$ thus
648 promoting its bioavailability to marine phytoplankton for primary production (31)
649 and against degradation by microbes (32). Mutation in genes that encode enzymes
650 for the biosynthesis of LPS in *E. coli* was shown to confer resistance to T7 phage
651 (33). Thus, organisms, even at the level of strains, seem to evolve to modify their
652 monosaccharide repertoire to increase fitness. In fact, selection pressure and
653 horizontal gene transfer events could be the reason for the monosaccharide
654 repertoire of bacteria far exceeding those of mammalian and other eukaryotes.
655
656 Genes encoding enzymes for the biosynthesis of Neu5Ac are found in 5% and
657 0.6% genomes of Alpha-proteobacteria and Actinobacteria, respectively; the
658 bacterial carbohydrate structure database had no Neu5Ac-containing glycan from
659 organisms belonging to this class/phylum (9). L-Rhamnose and L-fucose are found
660 in 16% of Delta- and Epsilon-proteobacteria genomes and in 25% of actinobacteria
661 genomes. However, very few L-rhamnose- and L-fucose-containing glycans from
662 these classes/phyla are deposited in the database leading to the inference that these
663 are rare sugars in this class/phylum. Thus, inferring monosaccharide usage based
664 on an analysis of experimentally characterized glycans can at best give a partial
665 picture.
666
667 Rare group monosaccharides are those which are found only in a few species, few
668 genera and few phyla. Reasons for acquiring Rare group sugars can at best be
669 speculative. For instance, Bac2Ac4Ac occurs at the reducing end of glycans N-
670 and O-linked to proteins (34) but the presence of Bac2Ac4Ac is not mandatory for
671 *C. jejuni* PglB, an oligosaccharyltransferase, since it can transfer glycans which
672 have Glc2NAc, Gal2NAc or Fuc2NAc also at the reducing end (35). Perhaps,
673 Bac2Ac4Ac provides resistance to enzymes like PNGase F that cleave off N-
674 glycans. L-rhamnose, Neu5Ac, L-Qui2NAc, Man2NAc and L-Ara4N are not used
675 by *Leptospira biflexa* (a non-pathogen) but are used by *Leptospira interrogans* (a
676 pathogen). It is tempting to infer that these monosaccharides impart virulence to
677 the latter but analysis of monosaccharides used by *E. coli* strains belonging to
678 multiple pathotypes (enterohemorragic, enteropathogenic, uropathogenic) did not

23

679  reveal any relationship between monosaccharides and their phenotype. Tyvelose,
680  paratose and abequose are 3,6-dideoxy sugars that belong to the Rare group. These
681  are found primarily in *Salmonella enterica*, *Yersinia pestis* and *Yersinia*
682  *pseudotuberculosis*. These are present in the O-antigen of *Y. pseudotuberculosis*
683  (36). *Y. pestis*, closely related to and derived from *Y. pseudotuberculosis*, lacks O-
684  antigen (rough phenotype) due to the silencing of O-antigen cluster (37). *Y.*
685  *enterocolitica*, also an enteric pathogen like *Y. pseudotuberculosis*, does not
686  contain these monosaccharides. Hence the role of these 3,6-dideoxy sugars in the
687  O-antigen of *Y. pseudotuberculosis* does not seem to be related to
688  enteropathogenicity.

690  Besides answering the question of the universality of glycan alphabet, this study
691  also has led to certain beneficial outcomes. L-rhamnose, mannose and L-
692  Pse5ac7Ac are found in *B. cereus, B. mycoides* and *B. thuringeinsis* but not in *B.*
693  *subtilis, B. amyloliquefaciens, B. licheniformis, B. velezensis and B. vallismortis.*
694  Such differences may be potentially be exploited towards taxonomic identification,
695  provided that these patterns hold true after analysis of a larger number of strains
696  from each of these species. Enzymes synthesizing monosaccharides that are
697  exclusive to a pathogen vis-à-vis its host can be identified as potential drug targets.
698  An illustrative example is of the non-hydrolyzing C2 epimerase: it mediates the
699  synthesis of UDP-Man2NAc, UDP-L-Qui2NAc, UDP-L-Fuc2NAc and UDP-
700  Man2NAc3NAc and is found in 60% of the archaeal+bacterial genomes but not in
701  humans (human genome was scanned for the presence of these pathways;
702  unpublished results). It has already been reported that inhibitors of this enzyme are
703  effective against methicillin-resistant *S. aureus* and a few other bacteria (38).
704  Based on the prevalence of this enzyme in all other phyla, inhibitors against this
705  enzyme would be promising broad spectrum antimicrobial therapies. As already
706  noted (39), knowledge of monosaccharide composition is also useful in ensuring
707  consistency of recombinant glycoprotein therapeutics. Knowledge of biosynthesis
708  pathways also allows cloning the entire cassette in a heterologous host for large
709  scale production of monosaccharides for commercial and research applications.

711  Thus, glycans show least evolutionary conservation among these three
712  macromolecules (40). Owing to their virtue of endowing distinction, existence of a

24

713  universal glycan alphabet is antithetical. Here, alphabet is used in the same sense
714  as its dictionary meaning, viz., a set of letters or symbols which combine to form
715  complex entities. In the case of glycans, structural diversity arises not only by the
716  set of monosaccharides an organism uses but also by linkage variations ($\alpha$1$\rightarrow$3,
717  $\beta$1$\rightarrow$4, etc.), branching and modifications (e.g., sulfation, acetylation, …).
718  Knowledge of the linkage types, branching patterns and modifications that an
719  organism uses will further our understanding of the biological roles of glycans.

**Author contributions**

PVB conceived and supervised the research, PS conceived the design and guided the development of GlycoPathDB. JS performed the research and developed the database. PVB and JS wrote the paper.

**Conflicts of interest**

The author(s) declare that there are no conflicts of interest

**Funding information**

This work received no specific grant from any funding agency

**Acknowledgments**

# References

738 
739  1.  Varki A. Biological roles of oligosaccharides: all of the theories are correct. Glycobiology.
740      1993;3(2):97–130.

741  2.  Mostowy RJ, Croucher NJ, De Maio N, Chewapreecha C, Salter SJ, Turner P, et al.
742      Pneumococcal Capsule Synthesis Locus cps as Evolutionary Hotspot with Potential to
743      Generate Novel Serotypes by Recombination. Mol Biol Evol. 2017 01;34(10):2537–54.

744  3.  Mostowy RJ, Holt KE. Diversity-Generating Machines: Genetics of Bacterial Sugar-
745      Coating. Trends Microbiol. 2018;26(12):1008–21.

746  4.  Gabius H-J, Roth J. An introduction to the sugar code. Histochem Cell Biol. 2017
747      Feb;147(2):111–7.

748  5.  Bravo D, Silva C, Carter JA, Hoare A, Alvarez SA, Blondel CJ, et al. Growth-phase
749      regulation of lipopolysaccharide O-antigen chain length influences serum resistance in
750      serovars of Salmonella. J Med Microbiol. 2008 Aug;57(Pt 8):938–46.

751  6.  Kalynych S, Morona R, Cygler M. Progress in understanding the assembly process of
752      bacterial O-antigen. FEMS Microbiol Rev. 2014 Sep;38(5):1048–65.

753  7.  Johannessen C, Koomey M, Børud B. Hypomorphic glycosyltransferase alleles and
754      recoding at contingency loci influence glycan microheterogeneity in the protein
755      glycosylation system of Neisseria species. J Bacteriol. 2012 Sep;194(18):5034–43.

756  8.  Kaltner H, Abad-Rodríguez J, Corfield AP, Kopitz J, Gabius H-J. The sugar code: letters
757      and vocabulary, writers, editors and readers and biosignificance of functional glycan-lectin
758      pairing. Biochem J. 2019 24;476(18):2623–55.

759  9.  Herget S, Toukach PV, Ranzinger R, Hull WE, Knirel YA, von der Lieth C-W. Statistical
760      analysis of the Bacterial Carbohydrate Structure Data Base (BCSDB): characteristics and
761      diversity of bacterial carbohydrates in comparison with mammalian glycans. BMC Struct
762      Biol. 2008 Aug 11;8:35.

763  10. Malmström A, Bartolini B, Thelin MA, Pacheco B, Maccarana M. Iduronic acid in
764      chondroitin/dermatan sulfate: biosynthesis and biological function. J Histochem Cytochem.
765      2012 Dec;60(12):916–25.

766  11. Tello M, Jakimowicz P, Errey JC, Freel Meyers CL, Walsh CT, Buttner MJ, et al.
767      Characterisation of Streptomyces spheroides NovW and revision of its functional
768      assignment to a dTDP-6-deoxy-D-xylo-4-hexulose 3-epimerase. Chem Commun (Camb).
769      2006 Mar 14;(10):1079–81.

770  12. Polizzi SJ, Walsh RM, Peeples WB, Lim J-M, Wells L, Wood ZA. Human UDP-α-D-
771      xylose synthase and Escherichia coli ArnA conserve a conformational shunt that controls
772      whether xylose or 4-keto-xylose is produced. Biochemistry. 2012 Nov 6;51(44):8844–55.

773    13.  Li Z, Mukherjee T, Bowler K, Namdari S, Snow Z, Prestridge S, et al. A four-gene operon
774         in Bacillus cereus produces two rare spore-decorating sugars. J Biol Chem. 2017 May
775         5;292(18):7636–50.

776    14.  Hwang S, Aronov A, Bar-Peled M. The Biosynthesis of UDP-D-QuiNAc in Bacillus cereus
777         ATCC 14579. PLoS ONE. 2015;10(7):e0133790.

778    15.  Ohno S. Evolution by Gene Duplication. Springer Science & Business Media; 2013. 171 p.

779    16.  Copley SD. Evolution of new enzymes by gene duplication and divergence. FEBS J. 2020
780         Apr;287(7):1262–83.

781    17.  Keinhörster D, George SE, Weidenmaier C, Wolz C. Function and regulation of
782         Staphylococcus aureus wall teichoic acids and capsular polysaccharides. Int J Med
783         Microbiol. 2019 Sep;309(6):151333.

784    18.  Pfannkuch L, Hurwitz R, Traulsen J, Sigulla J, Poeschke M, Matzner L, et al. ADP heptose,
785         a novel pathogen-associated molecular pattern identified in Helicobacter pylori. FASEB J.
786         2019;33(8):9087–99.

787    19.  Trachtenberg S. Mollicutes-wall-less bacteria with internal cytoskeletons. J Struct Biol.
788         1998 Dec 15;124(2–3):244–56.

789    20.  Khachane AN, Timmis KN, Martins dos Santos VAP. Dynamics of reductive genome
790         evolution in mitochondria and obligate intracellular microbes. Mol Biol Evol. 2007
791         Feb;24(2):449–56.

792    21.  Amano K, Tamura A, Ohashi N, Urakami H, Kaya S, Fukushi K. Deficiency of
793         peptidoglycan and lipopolysaccharide components in Rickettsia tsutsugamushi. Infect
794         Immun. 1987 Sep;55(9):2290–2.

795    22.  Rund S, Lindner B, Brade H, Holst O. Structural analysis of the lipopolysaccharide from
796         Chlamydia trachomatis serotype L2. J Biol Chem. 1999 Jun 11;274(24):16819–24.

797    23.  Peturova M, Vitiazeva V, Toman R. Structural features of the O-antigen of Rickettsia typhi,
798         the etiological agent of endemic typhus. Acta Virol. 2015 Sep;59(3):228–33.

799    24.  Theunissen C, Cnops L, Van Esbroeck M, Huits R, Bottieau E. Acute-phase diagnosis of
800         murine and scrub typhus in Belgian travelers by polymerase chain reaction: a case report.
801         BMC Infect Dis. 2017 13;17(1):273.

802    25.  Tamura A, Ohashi N, Urakami H, Miyamura S. Classification of Rickettsia tsutsugamushi
803         in a new genus, Orientia gen. nov., as Orientia tsutsugamushi comb. nov. Int J Syst
804         Bacteriol. 1995 Jul;45(3):589–91.

805    26.  Fuxelius H-H, Darby A, Min C-K, Cho N-H, Andersson SGE. The genomic and metabolic
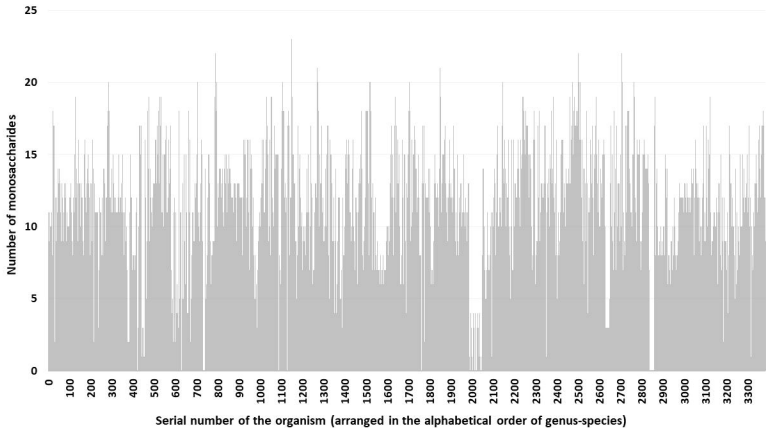806         diversity of Rickettsia. Res Microbiol. 2007 Dec;158(10):745–53.

807 27. Jones C. Revised structures for the capsular polysaccharides from Staphylococcus aureus
808      Types 5 and 8, components of novel glycoconjugate vaccines. Carbohydr Res. 2005 May
809      2;340(6):1097–106.

810 28. Moissl-Eichinger C, Pausan M, Taffner J, Berg G, Bang C, Schmitz RA. Archaea Are
811      Interactive Components of Complex Microbiomes. Trends Microbiol. 2018;26(1):70–85.

812 29. Lurie-Weinberger MN, Peeri M, Gophna U. Contribution of lateral gene transfer to the
813      gene repertoire of a gut-adapted methanogen. Genomics. 2012 Jan;99(1):52–8.

814 30. Lukácová M, Barák I, Kazár J. Role of structural variations of polysaccharide antigens in
815      the pathogenicity of Gram-negative bacteria. Clin Microbiol Infect. 2008 Mar;14(3):200–6.

816 31. Hassler CS, Schoemann V, Nichols CM, Butler ECV, Boyd PW. Saccharides enhance iron
817      bioavailability to Southern Ocean phytoplankton. Proc Natl Acad Sci USA. 2011 Jan
818      18;108(3):1076–81.

819 32. Zhang Z, Chen Y, Wang R, Cai R, Fu Y, Jiao N. The Fate of Marine Bacterial
820      Exopolysaccharide in Natural Marine Microbial Communities. PLoS ONE.
821      2015;10(11):e0142690.

822 33. Qimron U, Marintcheva B, Tabor S, Richardson CC. Genomewide screens for Escherichia
823      coli genes affecting growth of T7 bacteriophage. Proc Natl Acad Sci USA. 2006 Dec
824      12;103(50):19039–44.

825 34. Morrison MJ, Imperiali B. The renaissance of bacillosamine and its derivatives: pathway
826      characterization and implications in pathogenicity. Biochemistry. 2014 Feb 4;53(4):624–38.

827 35. Wacker M, Feldman MF, Callewaert N, Kowarik M, Clarke BR, Pohl NL, et al. Substrate
828      specificity of bacterial oligosaccharyltransferase suggests a common transfer mechanism
829      for the bacterial and eukaryotic systems. Proc Natl Acad Sci USA. 2006 May
830      2;103(18):7088–93.

831 36. Kenyon JJ, Cunneen MM, Reeves PR. Genetics and evolution of Yersinia
832      pseudotuberculosis O-specific polysaccharides: a novel pattern of O-antigen diversity.
833      FEMS Microbiol Rev. 2017 01;41(2):200–17.

834 37. Skurnik M, Peippo A, Ervelä E. Characterization of the O-antigen gene clusters of Yersinia
835      pseudotuberculosis and the cryptic O-antigen gene cluster of Yersinia pestis shows that the
836      plague bacillus is most closely related to and has evolved from Y. pseudotuberculosis
837      serotype O:1b. Mol Microbiol. 2000 Jul;37(2):316–30.

838 38. Xu Y, Brenning B, Clifford A, Vollmer D, Bearss J, Jones C, et al. Discovery of Novel
839      Putative Inhibitors of UDP-GlcNAc 2-Epimerase as Potent Antibacterial Agents. ACS Med
840      Chem Lett. 2013 Dec 12;4(12):1142–7.

841   39.   Mariño K, Bones J, Kattla JJ, Rudd PM. A systematic approach to protein glycosylation
842          analysis: a path through the maze. Nat Chem Biol. 2010 Oct;6(10):713–23.

843   40.   Varki A. Biological roles of glycans. Glycobiology. 2017;27(1):3–49.

844

## Common

**Simple sugars**
 Glucose
 Galactose
 Mannose
**C2-N-acetyl derivatives**
 Glc2NAc
 Gal2NAc
 Man2NAc
**Uronic acid derivatives**
 GlcA
 GalA
 Glc2NAcA
 Gal2NAcA
 Man2NAcA
**Deoxy derivative**
 L-Rhamnose

## Less common

**Simple sugars**
 Xylose
 L-Galactose
**Furanose form**
 Galactofuranose
**Deoxy derivative**
 L-Fucose
**C4-Amino derivative**
 L-Ara4N

## Rare

**Simple sugar**
 L-Galactose

**Amino / N-Acetyl derivatives**

| | | | |
|---|---|---|---|
| Qui2NAc | Fuc2NAc | Per | L-Qui2NAc |
| Qui3NAc | Fuc3NAc | Per4Ac | L-Fuc2NAc |
| Qui4NAc | Fuc4NAc | | L-Rha2NAc |
| Qui4NFo | Bac2Ac4Ac | L-Ara4NFo | |

**Uronic acid derivatives**

 ManA    Man2NAc3NAcA

**Deoxy derivatives**

| | | | |
|---|---|---|---|
| Rhamnose | Cillose | Fucofuranose | 6-Deoxytalose |
| Fucose | Cereose | Yelosamine | 6-Deoxygulose |
| | | | 6-Deoxy-L-talose |

**Dideoxy derivatives**

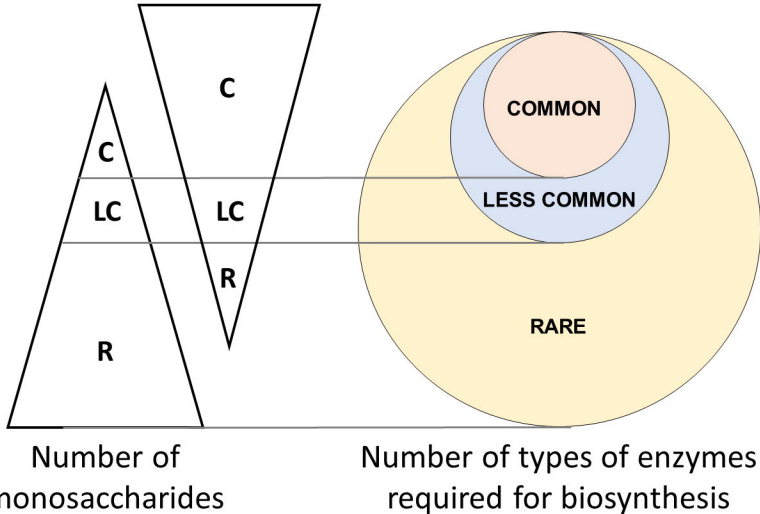 Paratose    Abequose    Tyvelose    L-colitose    L-Ascarylose

**Heptoses**

 L-*glycero*-β-D-*manno*-Heptose         D-*glycero*-α-D-*manno*-Heptose
 6-Deoxy-α-D-*manno*-heptose             6-Deoxy-α-D-*altro*-heptose

**9-Carbon sugars**

 Neu5Ac    Leg5Ac7Ac    L-Pse5Ac7Ac

Prevalence of monosaccharides

C

C

LC

LC

R

R

Number of monosaccharides

COMMON

LESS COMMON
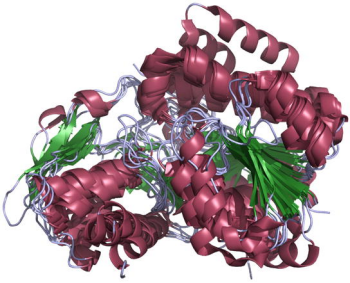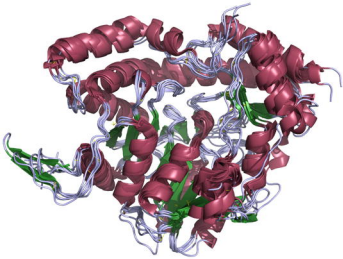
RARE

Number of types of enzymes required for biosynthesis

Glutamine amidotransferases
Nucleotidyltransferases
Cupin superfamily
SDR superfamily
α-D-phosphohexomutase superfamily
Non-hydrolysing C2-epimerases
N-acetyltransferases (LβH-domain containing)
Type II/III phosphomannose isomerase

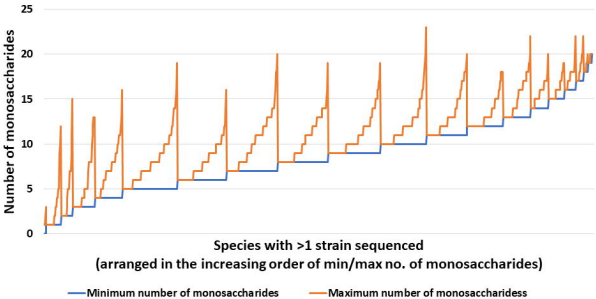PLP-dependent enzymes
UDP-galactose mutase family

Deaminases
N-formyltransferases
N-acetyltransferases (GNAT family)
Kinase
SAM-dependent methyltransferases
Hydrolases
PEP-utilizing synthetases
Hydrolyzing C2-epimerases
HAD superfamily
Sedoheptulose 7-phosphate isomerases
Gfo/Idh/MocA family

Types of enzymes

**Y-axis:** %age of genomes in each phyla

**Y-axis scale:** 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100

**X-axis categories:** Man2NAcA, Galactofuranose, L-Ara4N, L-Ara4NFo, L-Arabinose, L-Fuc2NAc, L-Qui2NAc, L-Rha2NAc, Man2NAc3NAcA, Qui2NAc, Xylose, Bac2Ac4Ac, Fuc2NAc, Yelosamine, 6-Deoxy-L-talose, Fuc3NAc, Fuc4NAc, Fucofuranose, Fucose, Qui3NAc, Qui4NAc, Qui4NFo, L-Glycero-b-D-manno-Hep, 6-Deoxy-gulose, Cereose/Clllose, L-Ascarylose, Paratose/Abequose, Tyvelose, Leg5Ac7Ac, L-Pse5Ac7Ac, Neu5Ac, 6-Deoxy-a-D-altro-Hep, 6-Deoxy-a-D-manno-Hep, Rhamnose/6-deoxy-talose, D-Glycero-a-D-manno-Hep, L-Fucose, L-Galactose, ManA, Per, Per4Ac, L-Colitose

**Legend:**
- Nitrospirae (9)
- Acidobacteria (10)
- Aquificae (15)
- Thermotogae (33)
- Fusobacteria (50)
- TACK group archaea (92)
- Spirochaetes (95)
- PVC group (200)
- Euryarchaeota (208)
- FCB group (497)
- Proteobacteria (6985)
- Terrabacteria (4707)