1

# Phenotypic and genomic diversification in complex carbohydrate degrading human gut bacteria

Nicholas A. Pudlo[1*], Karthik Urs[1*], Ryan Crawford[2], Ali Pirani[1], Todd Atherly[3,4], Roberto Jimenez[5], Nicolas Terrapon[6,7], Bernard Henrissat[6,7,8], Daniel Peterson[5,9], Cherie Ziemer[3,4], Evan Snitkin[1] and Eric C. Martens[1]

[1]Department of Microbiology and Immunology, University of Michigan Medical School, Ann Arbor, MI 48109

[2]Department of Computational Medicine and Bioinformatics, University of Michigan Medical School, Ann Arbor, MI 48109

[3]Iowa State University, Dept. of Animal Science, Ames, IA

[4]United States Department of Agriculture Agricultural Research Station, Ames, IA

[5]University of Nebraska, Department of Food Sciences, Lincoln, NE

[6]Aix Marseille Univ, CNRS, UMR7257 AFMB, Marseille, France

[7]INRAE, USC1408 AFMB, Marseille, France

[8]Department of Biological Sciences, King Abdulaziz University, Jeddah, Saudi Arabia

[9]Johns Hopkins University, Department of Pathology, Baltimore, MD

*Authors contributed equally to this work

Correspondence to: emartens@umich.edu

**Running Title:** *Bacteroidetes* carbohydrate utilization

40    **Abstract**

41        Symbiotic bacteria are responsible for the majority of complex carbohydrate digestion in

42    the human colon. Since the identities and amounts of dietary polysaccharides directly impact the

43    gut microbiota, determining which microorganisms consume specific nutrients is central to

44    defining the relationship between diet and gut microbial ecology. Using a custom phenotyping

45    array, we determined carbohydrate utilization profiles for 354 members of the Bacteroidetes, a

46    dominant saccharolytic phylum. There was wide variation in the numbers and types of substrates

47    degraded by individual bacteria, but phenotype-based clustering grouped members of the same

48    species indicating that each species performs characteristic roles. The ability to utilize dietary

49    polysaccharides and endogenous mucin glycans was negatively correlated, suggesting exclusion

50    between these niches. By analyzing related *Bacteroides ovatus/xylanisolvens* strains that vary in

51    their ability to utilize mucin glycans, we addressed whether gene clusters that confer this

52    complex, multi-locus trait are being gained or lost in individual strains. Pangenome

53    reconstruction of these strains revealed a remarkably mosaic architecture in which genes

54    involved in polysaccharide metabolism are highly variable and bioinformatics data provide

55    evidence of interspecies gene transfer that might explain this genomic heterogeneity. Global

56    transcriptomic analyses suggest that the ability to utilize mucin has been lost in some lineages of

57    *B. ovatus* and *B. xylanisolvens*, which still harbor residual gene clusters that are involved in

58    mucin utilization by strains that still actively express this phenotype. Our data provide insight

59    into the breadth and complexity of carbohydrate metabolism in the microbiome and the

60    underlying genomic events that shape these behaviors.

61

62

63

64   **Introduction**

65       Microbial communities in the distal intestines of humans and other mammals play critical

66   roles in the digestion of dietary polysaccharides (1-3). Unlike proteins, lipids and simple sugars,

67   which can be assimilated in the small intestine, the vast majority of non-starch polysaccharides

68   (fibers) transit undegraded to the distal gut due to a lack of requisite enzymes encoded in the

69   human genome (4). Microbial transformation of dietary fiber polysaccharides into host-

70   absorbable organic and short chain fatty acids is a beneficial process that unlocks otherwise

71   unusable calories from our diet (5), shapes the composition and behavior of the gut microbial

72   community (6-8), provides preferred nutrients directly to the colonic epithelium (9-11) and

73   shapes the development of immune cell populations (12, 13).

74       The abundance of dietary fiber in the mammalian diet, and the substantial chemical

75   diversity within this class of molecules, provides a prominent selective pressure that drives

76   genome evolution and diversification within symbiotic bacterial populations. The genomes of

77   individual human gut bacteria frequently encode dozens-hundreds more polysaccharide-

78   degrading enzymes than human secrete into the gastrointestinal tract, reflecting gut microbial

79   adaptations to degrade dietary fibers (3, 4). As examples, the genomes of a few well-studied

80   Gram-negative *Bacteroides* (*B. thetaiotaomicron*, *B. ovatus* and *B. cellulosyliticus*) encode

81   between 250 and over 400 CAZymes that collectively equip them to target nearly all commonly

82   available dietary polysaccharides (14-16). However, none of these three species is by itself

83   capable of degrading all available polysaccharides, a conclusion that was supported by early

84   phenotypic surveys of cultured human gut bacteria that encompassed species from other phyla

85   (17, 18). These findings suggest that individual microbes fill multiple, specific carbohydrate

86   degradation niches and that a diverse community is required to ensure degradation of the entire

87    repertoire of dietary fibers. Given that hundreds of different microbial species typically coexist in

88    an individual over long time periods (19), it is important to understand how many different

89    polysaccharide metabolism pathways are present within the individual microbial species that

90    compose a community and how these traits are represented across strains and species. If some

91    species possess very similar phenotypic abilities, they may be functional surrogates or compete

92    for similar niches and therefore seldom co-occur.

93        Members of the Bacteroidetes are often among the most numerous bacteria in the colonic

94    microbiota of people from industrialized countries (19-21). These bacteria are well appreciated

95    for their abilities to degrade a broad range of polysaccharides (16-18, 22, 23) and modify disease

96    states in a bacterial species-specific fashion (24-26). In this study, we empirically measured the

97    abilities of members of 29 different Bacteroidetes species to grow on a custom panel of

98    carbohydrates that span the diversity of plant, animal and microbial polysaccharides. Our results

99    reveal a wide range of metabolic breadth between different species, indicating that some have

100   evolved to be carbohydrate generalists while others have become metabolically specialized to

101   target just one or a few nutrients. Pangenome analysis of several related strains provides insight

102   into the evolutionary events that shape carbohydrate utilization among these important symbionts

103   and reveals a dizzying mosaic of heterogeneity at the level of discrete gene clusters mediating

104   polysaccharide metabolism. Based on analysis of several variable loci, we provide evidence to

105   support a mechanism of lateral gene transfer that may account for this mosaic architecture. Our

106   results provide a glimpse into the metabolic breadth and diversity of an important group of

107   human gut bacteria towards polysaccharide metabolism. Given the large amount of genomic and

108   metagenomic sequence information that has been generated from the human microbiome,

109     phenotypic studies such as the one presented here represent important next steps in deciphering

110     the functionality of these organisms in their native gut habitat.

111

112     **Results**

113         Phenotypes are the ultimate measures of biological function. However, large-scale

114     phenotypic analyses are still uncommon in surveys of the human gut microbiome, which have

115     instead relied on sequence-based approaches to infer function, often with substantial uncertainty.

116     This lack of phenotypic information is partly due to a lack of high-density (*e.g.*, strain level)

117     culture representation for the dominant taxa combined with a lack of defined growth conditions

118     to measure behavior of these organisms. With the resurgence of gut microbial culturing, both of

119     these gaps have begun to close (27-30), revealing an urgent need for scalable platforms to define

120     the actual behavior of these organisms. To address this gap, we assembled a collection of human

121     and animal gut Bacteroidetes and constructed a custom anaerobic phenotyping platform centered

122     around carbohydrate metabolism, a key function that symbiotic gut microorganisms contribute to

123     mammalian digestion (4). This array consists of 45 different carbohydrates (30 polysaccharides

124     and 15 monosaccharides) that span the repertoire of common sugars and linkages present in

125     dietary plants and meat, as well as host mucosal secretions and some rare nutrients consumed in

126     regional populations or as food additives (see **Fig. S1** for a summary of polysaccharide

127     structures).

128         The carbohydrate utilization abilities of 354 different human and animal Bacteroidetes

129     strains were measured by individually inoculating each into this custom growth array and

130     automatically monitoring anaerobic growth every 10-20 min for four days (see *Materials and*

131    *Methods*). Based on 16S rRNA gene sequence for each strain, this collection encompasses 29

132    different species based on the requirement that each strain possesses ≥98% 16S rRNA gene

133    identity to a named type strain in a given species (**Table S1,** note that all but three strains, which

134    were all related to each other and to *Bacteroides uniformis*, met this criterion). The resulting

135    31,860 individual growth curves were first inspected manually and then subjected to automated

136    analysis to quantify total growth and growth rate parameters for each substrate (see *Materials*

137    *and Methods*). A normalization scheme was employed to compensate for general growth

138    differences in the two different defined medium formulations employed (see **Table S1** for a full

139    list of strains assayed and all raw and normalized growth measurements, **Fig. S2** for analysis of

140    replicates).

141    **Members of the same species possess similar carbohydrate utilization profiles**

142         Growth results are summarized in **Figs. 1**, **2** and **S3**. Whether considered from the

143    perspective of how many species degrade a particular polysaccharide (**Fig. 1A**), or how many

144    individual polysaccharides are targeted by members of a particular species (**Fig. 1B**), there was

145    substantial variability in carbohydrate utilization among the organisms surveyed (range, 1-28

146    polysaccharides degraded per strain; mean, 15.6). Some polysaccharides like soluble

147    starch/glycogen were degraded by a majority of the species tested, yet others like the edible

148    seaweed polysaccharides carrageenan and porphyran were used by just one or two strains.

149         Given the diversity in observed carbohydrate utilization phenotypes, we wished to

150    address if closely related strains display similar abilities or instead if strains of the same species

151    have diverged from one another. To assist in visualizing the overall trends in carbohydrate

152    utilization across this phylum, we performed unsupervised clustering of the strains based on their

153    carbohydrate utilization profiles. While many species are not deeply represented by multiple

154   strains, clustering based on a combination of normalized growth and rate measurements largely

155   grouped strains of the same species together (**Fig. 2**) and, as expected, this was driven mostly by

156   polysaccharide utilization abilities (**Fig. S4**).

157   Our data reveal that strains belonging to several individual species possess more similar

158   polysaccharide degrading abilities to each other compared to their more distant relatives, a

159   finding that has importance for interpreting or predicting function based on community

160   sequencing data. As examples, all 56 strains of *B. fragilis* clustered together, reflecting their

161   generally restricted abilities to utilize forms of soluble starch/glycogen, inulin and mucus *O*-

162   glycans. Likewise, all 36 strains of *B. uniformis*, a species with broader metabolic capacity that

163   includes digestion of plant cell wall hemicelluloses, were also grouped together into a single

164   branch. The inclusivity of these groupings was generally independent of the time period when

165   strains were isolated or whether they were isolated from humans or other mammals (**Fig. 2**).

166   Another important feature of the observed species clustering is that the grouping does not

167   mirror the overall phylogeny of the gut Bacteroidetes. Rather, phylogenetically separated species

168   often group adjacent to one another based on similarities in carbohydrate metabolism (*e.g.*, *B.*

169   *ovatus/xylanisolvens* and *B. cellulosilyticus;* and *B. vulgatus/dorei* and *B. fragilis*; see **Fig. 3A**

170   for a phylogenetic tree based on conserved housekeeping genes) (31, 32). In the latter case, it is

171   interesting to directly compare *B. fragilis* and *B. vulgatus/dorei*, two groups with deep strain

172   representation (**Fig. 2**). Despite being phylogenetically more distant, these species possess very

173   similar phenotypic patterns that reflect degradation of soluble starch and similar molecules

174   (glycogen, pullulan), inulin and mucin *O*-glycans. The major distinguishing feature between

175   these groups is the presence of some, often-weak, pectin utilization among strains of *B.*

176   *vulgatus/dorei*.

177    Some polysaccharides, especially those present in the cell walls of dietary plants, occur in

178    the same physical context and presumably traverse the gut together, potentially exerting selective

179    pressure for bacteria to use them simultaneously. To test for co-occurrence of traits, we

180    performed a pairwise correlation analysis to determine the extent to which any two

181    polysaccharides were co-utilized by the same strain (**Fig. S5**). The presence of two different

182    soluble starches (potato and maize amylopectin) and two starch-like glycans (glycogen and

183    pullulan) provide an internal control since they are essentially identical in their sugar and linkage

184    chemistry but vary in the proportion and placement of branches as well as polymer length,

185    crystallinity and solubility (**Fig. S1**). These four molecules are utilized through a single

186    degradation/transport system in the type strain of *B. thetaiotaomicron*, which was included in our

187    study (33). As expected, the abilities to use these four polysaccharides were among the strongest

188    positive correlations (between 44-75%); although, there was not a perfect correlation suggesting

189    that some finer adaptation may exist even for different structural forms of a chemically similar

190    molecule.

191    We also observed positive correlations in the ability to use components of two different

192    groups of plant cell wall polysaccharides (pectins and hemicelluloses), as well as animal tissue

193    glycosaminoglycans, despite the fact that the polysaccharides within each of these groups often

194    possess different chemical structures (**Fig. S1**). In the case of the hemicelluloses, there was even

195    some apparent separation based on dicotyledonous vs. monocotyledonous sources. The

196    predominantly dicot hemicelluloses (**Fig. 2**, blue labels) and monocot hemicelluloses (**Fig. 2**,

197    green labels) show some exclusivity with respect to the bacteria that utilize them. Many *B.*

198    *ovatus/B. xylanisolvens* strains lack the ability to utilize the three dicot hemicelluloses (GalM,

199    GlcM, XyG); whereas the ability to degrade those from monocots (OSX, WAX, BBG) is more

200     evenly distributed. *B. uniformis* has a partially opposite pattern, preferring substrates from dicots,

201     while only degrading one of the two major monocot structures (BBG) and poorly degrading the

202     two xylans tested (OSX, WAX). Similar observations were also made for pectins and GAGs and

203     could reflect adaptations to simultaneously harvest different nutrients from digesta particles

204     derived from dicot plant cell walls or animal tissue ingested in a carnivorous diet.

205

206     **Specialization for mucus *O*-linked glycans**

207          The most noteworthy correlation between polysaccharide utilization traits was observed

208     between utilization of host-produced mucin *O*-glycans and many of the other polysaccharides

209     tested. Growth on a total of 19/30 polysaccharides showed negative correlations with the ability

210     to utilize *O*-glycans, with the strongest negative correlations being between *O*-glycans and the

211     seven different hemicelluloses (**Fig. S5**). This negative correlation is easily observed by

212     comparing the rightmost column in **Fig. 2** (*O*-glycan utilization) with the respective columns for

213     hemicellulose degradation. Because this trend was observed across several species, it suggests

214     that there could be a more general exclusive relationship between the two niches associated with

215     foraging on mucus and hemicellulose. This idea is further supported by experiments described

216     below, which suggest that isolates of *B. ovatus* and *B. xylanisolvens*, both adept hemicellulose

217     consumers, are in the process of losing the ability to degrade *O*-glycans, relative to an ancestor

218     that contained multiple gene clusters involved in the metabolism of these structures.

219          Interestingly, the mucin *O*-glycan mixture was the only substrate for which we observed

220     absolute metabolic specialization among the substrates tested. A single, and only available strain

221     of *Barnesiella intestinihominis* exhibited the ability to exclusively utilize mucin *O*-glycans, along

222     with a subset of the sugars that are contained in these structures (**Table S1**). Three strains of

223    *Bacteroides massiliensis* exhibited similar behavior with very strong growth on mucin *O*-glycans

224    and only weak growth on soluble starches and a few other polysaccharides (**Table S1**). These

225    three *B. massiliensis* strains were also restricted in the repertoire of simple sugars they could

226    metabolize with this list being limited to those found in mucin and other host glycans (galactose,

227    *N*-acetylgalactosamine, *N*-acetylglucosamine, *N*-acetylneuraminic acid and L-fucose; weak

228    fructose utilization by one strain was the only exception). Members of these two species are

229    poorly represented in culture collections and remain lightly studied. However, their specific

230    adaptations for host mucin glycans may render them important members of the microbiota,

231    potentially thriving at the interface between the gut lumen and host tissue and relying exclusively

232    on the host to be sustained. The continuous supply of mucin *in vivo* could explain why some

233    species have become specialized for it as a nutrient, whereas dietary fiber degraders may need to

234    be more generalist since the substrates available to them change with the host's meals.

235

236    **Pangenome reconstruction reveals extensive genetic diversification among related**

237    ***Bacteroides***

238          With a view of the carbohydrate utilization traits present in our gut Bacteroidetes

239    collection, we next sought to determine if certain variable traits were being gained or lost within

240    strains of certain species and if available genomes provide insight into the mechanisms driving

241    genomic adaptations to particular nutrients. Connections between polysaccharide utilization

242    phenotypes and the underlying genes involved have been systematically explored for a few

243    *Bacteroides* species (*B. thetaiotaomicron, B. ovatus* and *B. cellulosilyticus*) with partial analyses

244    in others (6, 16, 22, 23, 34-37). These studies have revealed that, in essentially all cases, the

245    ability to degrade a particular polysaccharide is conferred by one or more clusters of co-

246     expressed genes termed polysaccharide utilization loci (PULs) (38). PULs share defining

247     features such as genes encoding homologs of outer membrane TonB-dependent transporters

248     (SusC-like), surface glycan-binding proteins (SGBPs; or SusD- and SusE/F-like), usually an

249     associated sensor/transcriptional regulator and one or more degradative CAZymes (glycoside

250     hydrolase, GH; polysaccharide lyase, PL; carbohydrate esterase, CE), as well as other enzymes

251     like sulfatases or proteases. Since the presence of one or more cognate PULs is required to utilize

252     a given polysaccharide and these genes typically exhibit large increases in gene expression in

253     response to their growth substrate, we rationalized that we could focus on traits that were

254     variable in closely related strains and locate the associated PULs by transcriptomic analysis to

255     gain insight into the basis of their acquisition or loss.

256            To test this, we focused on members of two closely related species, *B. ovatus* (*Bo*) and *B.*

257     *xylanisolvens* (*Bx*), for which there is noticeable inter-strain variation in the ability to use mucin

258     *O*-glycans (**Figs. 2, 3**). Investigation of these two species also benefits from substantial culture

259     depth and many strains with available sequences. The *O*-glycans attached to mucins represent a

260     diverse family of over one hundred different structures (39), albeit with common linkage patterns

261     (**Fig. S1**). Correspondingly, the ability to utilize these glycans is a complex trait, involving

262     simultaneous expression of at least 6-13 different *O*-glycan inducible PULs in *B.*

263     *thetaiotaomicron, B. massiliensis*, *B. fragilis* and *B. caccae* (6, 22, 35). Quantification of *O*-

264     glycan growth for individual *Bo* and *Bx* strains was widely variable (**Fig. 3B**). One hypothesis to

265     explain this variability is that some *Bo* and *Bx* strains have gained the ability to utilize *O*-glycans

266     relative to an ancestor that lacked this phenotype. If so, the PULs they express during *O*-glycan

267     degradation might be unique to their genomes and may indicate lateral gene transfer (LGT) as

268     has been the case for acquisition of phenotypes such as porphyran, agarose and λ-carrageenan

269    utilization in gut *Bacteroides*, which are all components of integrative conjugative elements or

270    mobilizable plasmids (31, 40). An alternative hypothesis is that some *Bo* and *Bx* strains are in the

271    process of losing this ability from a common ancestor. If so, the genomes of non-degraders may

272    still contain some PULs that are homologous to those present in more proficient *O*-glycan-

273    degrading strains, but these strains may have lost a key step(s) that has eroded their ability to

274    express this phenotype.

275        To distinguish these hypotheses, we selected seven strains (black arrows in **Fig. 3B**) that

276    vary in their ability to degrade *O*-glycans and for which genome sequences exist. Note that three

277    strains that degrade *O*-glycans were initially chosen because they were among the strongest

278    degraders in our dataset with sequenced genomes when we initiated these experiments. We later

279    identified strains with better *O*-glycan growth abilities and address one of these (strain H59)

280    separately below. Four of the selected strains were *Bo* (two positive and two negative for *O*-

281    glycan degradation); three strains were *Bx* (one weakly positive and two negative for *O*-glycan

282    degradation). One of these strains (*B. xylanisolvens* XB1A) has a finished circular genome and

283    was used as a scaffold to align the remaining six draft genome sequences, with manual curation

284    (see *Material and Methods*), resulting in a nearly contiguous pangenome sequence that captures

285    the spatial arrangement of homologous and variable genes that are present in these seven strains

286    (**Table S2**, **Fig. S6**).

287        Analysis of the *Bo/Bx* pangenome revealed remarkable variability in gene content among

288    just the seven strains used. A total of 12,960 different genes were delineated based on $\geq$90%

289    identity in their translated amino acid sequence (**Table S2**). Remarkably, only 2,264 (17.5%) of

290    these genes were shared among all seven strains. The largest proportion of genes (7,244; 55.9%)

291    was only present in one of the seven strains. Separating two major classes of core PUL functions,

292    SusC/D homologs and degradative CAZymes (GH, PL and CE), revealed that these key

293    components of Bacteroidetes polysaccharide metabolism were also heavily represented in the

294    "accessory gene" pool that is not common to all strains (**Fig. 4A**).

295         Through informatics-based and manual annotation of gene clusters containing typical

296    PUL functions, we delineated between 180-236 different PULs in the reconstructed pangenome

297    (ambiguity is caused by many PULs occurring adjacent to each other; although in many cases

298    separation of adjacent PULs according to individual genomes allowed us to make more precise

299    delineations, **Table S3**). Direct comparison of the *O*-glycan-degrading and non-degrading strains

300    revealed that there was a substantial number of genes (3,351) that were unique to the three *O*-

301    glycan degrading strains, including genes belonging to 51 PULs (**Fig. 4B**). However, such a

302    distribution in gene content might be expected given the overall large proportion of non-core

303    genes in these seven strains and there was correspondingly no indication that all three *O*-glycan-

304    degrading strains shared overlapping PULs with each other: no PULs were common to all three

305    *O*-glycan degraders and only five PULs were shared by any two strains (**Fig. 4C**). Considering

306    that there are 51 total PULs that are unique to the mucin-degrading strains, if these strains have

307    gained the ability to degrade *O*-glycans from an ancestral lineage that lacked this ability it likely

308    occurred by acquisition of separate gene clusters. To more directly distinguish between the two

309    hypotheses given above, we performed transcriptional profiling on all three *O*-glycan degrading

310    strains to determine if the PUL genes that they express during *O*-glycan degradation are indeed

311    unique to these strains.

312         Compared to reference growth in minimal medium containing glucose (MM-glucose), the

313    *Bx* D22, *Bo* 3-1-23 and *Bo* D2 strains activated expression of 196, 227 and 359 total genes more

314    than 10-fold and these gene lists included components of 14, 19, and 42 different PULs,

315 respectively (**Tables S4-6**). As expected from studies in other *Bacteroides*, these PULs were

316 scattered throughout the genome (**Fig. S7**), suggesting that they are autonomously regulated in

317 response to glycan cues present in the *O*-glycan mixture. Strikingly, the majority of PULs that

318 contained *O*-glycan-activated genes (63/75, 84%) were *not* unique to the *O*-glycan degrading

319 strains (**Tables S4-S6, Fig. S7**). Moreover, in each of the three strains analyzed, the most highly

320 upregulated PULs were also often shared with non-mucin degrading strains. These observations

321 lend support to the hypothesis that strains of *Bo* and *Bx* are in the process of losing the ability to

322 utilize *O*-glycans relative to a common ancestor that possessed a more expansive gene repertoire

323 to successfully access these nutrients. However, we cannot rule out that individual non-degrading

324 strains are separately acquiring PULs that are associated with mucin degradation and retaining

325 them without the full benefit that presumably occurs with the ability to fully execute this growth

326 phenotype. This latter idea is consistent with inter-species PUL exchange observations

327 elaborated on below.

328 Finally, because we subsequently identified a *B. ovatus* strain (NLAE-zl-H59, red arrow

329 in **Fig. 3B.**) with a substantially higher ability to use *O*-glycans relative to the strains used for

330 pangenome construction, we performed additional RNA-seq analysis on this strain. Compared to

331 a glucose reference, this strain activated 373 total genes in response to *O*-glycans, including

332 genes from 30 different PULs (**Table S7**). Among these, 26 activated PULs were also present in

333 one of the seven strains in our pangenome and 24 were homologous to PULs in strains that did

334 not degrade *O*-glycans. However, this strain did activate expression of genes within four PULs

335 that were completely unique to its genome compared to the seven strains used for pangenome

336 reconstruction, suggesting that it could possess additional genes that augment its ability to grow

337 on mucin *O*-glycans. This increased PUL expression could be responsible for the enhanced

338    growth of the H59 strain on *O*-glycans, especially if genes included within these unique PULs

339    are responsible for key metabolic steps required for efficient *O*-glycan utilization.

340

341    **Evidence that intergenomic recombination has driven *Bacteroides* pangenome evolution**

342            Similar to other bacteria, we observed that many accessory genes in the *Bo* and *Bx*

343    pangenome are located in contiguous clusters or "islands" often involving PULs or capsular

344    polysaccharide synthesis gene cluster (**Fig. S6**, **Table S2**). In contrast to previously identified

345    *Bacteroides* PULs that have more obviously been subjects of lateral transfer (31, 40, 41) and are

346    associated with integrative and conjugative elements (ICEs), most of the variable genomic

347    regions that we identified were not associated with functions indicative of mobile DNA. Instead,

348    these regions are often precisely located in between one or more core genes (*i.e.*, those common

349    to all seven strains; herein referred to as "genomic nodes") that flank each side of the variable

350    gene segment (**Fig. 5A,B**).

351            Several intergenomic transfer mechanisms might cause the observed mosaic structure of

352    the *Bo-Bx* pangenome. The first possibility is movement of variable genes into a recipient

353    genome by direct conjugation of individual, mobile ICEs. While such events would be expected

354    to leave behind residual genes involved in mobilization and transfer, which were not observed,

355    these DNA vehicles are known to target a subset of core genes, such as tRNAs (41), and may

356    have undergone subsequent genomic deletion events that eliminated the mobile DNA. Two other

357    known mechanisms of bacterial LGT are natural competence and phage transduction, neither of

358    which has been observed in members of the Bacteroidetes.

359            A final potential mechanism is direct conjugation of the chromosome from a donor

360    bacterium into a related recipient, followed by subsequent homologous recombination between

361    flanking nodes to add or delete intervening DNA in the recipient genome (**Fig. 5C**). A

362    mechanism that is conceptually similar to high-frequency recombination (Hfr) transfer in *E. coli*

363    has been described for *B. thetaiotaomicron* and *B. fragilis* and involves chromosomal ICEs,

364    which may have lost their ability to transfer autonomously by circularizing from the genome and

365    instead act as transfer initiation points to conjugate a donor genome into a recipient (42-44). If

366    such a mechanism was more broadly active in LGT between *Bacteroides*, we would expect that

367    some of the core/node genes involved would reflect sequence identities that were more similar to

368    the donor and this difference would be more easily detectable if the transfer was between

369    members of different species like *Bo* and *Bx*. Moreover, such transfer events could either result

370    in introduction of new genes into the recipient or elimination of genes depending on the genetic

371    content in between recombination nodes from the donor chromosome.

372         To test this hypothesis, we took a bioinformatics approach aimed at first identifying high-

373    confidence examples of inter-species recombination involving core genes and then assessed

374    whether those genes were associated with co-transfer of adjacent accessory genes (**Fig. 6A**). We

375    collected a dataset of 33 *Bo* and *Bx* genomes, which represent a subsample of the isolates for

376    which we generated phenotypic data. We identified a set of 1,384 core genes—expectedly

377    smaller than the core genome of the seven strains used above due to additional strains being

378    added—that are present as a single copy in all members of both species. To identify cases of

379    putative inter-species LGT via homologous recombination at core genes, we searched for

380    instances in which a core gene sequence was more similar to the corresponding sequence from

381    the other species than to sequences of the species to which a strain belonged. To this end, for

382    each allele of each core gene, we calculated the median distance to all other alleles of both

383    species (**Fig. 6B**). We identified instances where the median distance to the same species was

384     high and median distance to the opposite species was low and used these genes as markers for

385     putative LGT core loci. To identify additional accessory genes that may have been

386     simultaneously transferred, we searched for instances in which genes were perfectly syntenic and

387     collinear between each genome with a putative LGT core gene and genomes of the opposite

388     species. Among these candidate LGT loci, we then investigated if any of these transfer events

389     have resulted in pan-genome diversification, which we defined as the presence of any accessory

390     gene(s) that was only observed adjacent to a core gene with evidence of LGT based on the above

391     criteria.

392         In total, we identified 29 different loci at which exchange of core genes appeared to have

393     occurred and LGT accessory genes were identified, including seven that appeared to involve

394     transfer of PULs (**Fig. 6C**, **Fig. S8**). Similar numbers of potentially transferred loci were

395     identified for each species, with 16 loci in *Bx* and 13 loci in *Bo*. Within the identified HGT

396     events, variable numbers of HGT accessory genes were found within the loci ranging from one

397     to thirteen genes (**Fig. 6C**, **Fig. S8**). More genes (57 total) appeared to be transferred into *Bo*

398     than into *Bx* (36 total).

399         Finally, we determined if any of the identified LGT events could explain differential

400     phenotypes measured by our high throughput growth assay by modifying the complement of

401     PULs in individual genomes. As a specific example, we focused on a PUL that was previously

402     associated with β-mannan degradation (23, 45) that was among our candidate loci with evidence

403     of transfer from a *Bx* ancestor into two *Bo* strains. The presence of this PUL (PUL-A in **Fig. 6C,**

404     **Fig. S9A**) was observed in all strains with the ability to grow on the β-mannan galactomannan

405     (GalM), including two strains of *Bo* (ATCC8483 and CL02T12C04) for which the flanking node

406     regions were more similar to *Bx*. We previously showed that deletion of this PUL from *B. ovatus*

407    ATCC8483 eliminated growth on GalM and glucomannan (GluM) (45), suggesting that it was

408    both acquired from a *Bx* strain and conferred growth on these two β-mannans. However, the

409    presence of this PUL was not perfectly correlated with growth on GalM and several strains that

410    lacked PUL-A still exhibited robust growth. Thus, we searched for other PULs that harbor GH26

411    family enzymes and determined that all of the other strains that grow on GalM, but lack PUL-A,

412    harbor another candidate GalM PUL (PUL-B, **Fig. S9A**) at a different genomic location and

413    some strains possess both (**Fig. 6C**). Gene expression analysis by qPCR revealed that PUL-B

414    was highly expressed in strains that lacked PUL-A during growth in GalM (**Fig. S9B**) and every

415    strain that grew robustly on GalM had at least one of these two PULs. While we had previously

416    shown that PUL-A was required for GluM growth in *B. ovatus* ATCC8483, there were a number

417    of other strains (red "+" symbols in **Fig. 6C**) that displayed some ability to grow only on GlcM,

418    while lacking both of the GalM-associated PULs, suggesting the presence of additional, partially

419    orthogonal PULs that confer the ability to grow on variant β-mannans. Such a presence of

420    multiple orthologous PULs that confer the same or similar functions, and some which may be

421    moving between genomes of related species by the putative LGT mechanisms noted above,

422    complicates the process of understanding the genotype-phenotype relationships in human gut

423    Bacteroidetes, but will need to be resolved to make better functional predictions from sequence-

424    based data.

425

426    **Discussion**

427    In this study we leveraged a scalable, high-throughput quantitative growth platform to

428    characterize the phenotypic abilities that are present in a sample of hundreds of Bacteroidetes

429    strains from the human and animal gut. Our anaerobic screening technique is directly applicable

430    to other bacterial phyla from the human gut and other environments. Moreover, it can be adapted

431    to include new polysaccharides or to focus on different nutrient utilization or chemical resistance

432    phenotypes. The current study, in concert with future applications of phenotypic screening, will

433    help close the gap between our largely sequence-based view of the human gut microbiota and the

434    functions that its members provide. However, instances like the ones investigated here for mucin

435    glycan and β-mannan utilization by *Bacteroides* serve as a warning that the presence or absence

436    of genes that are experimentally associated with a particular function do not always indicate that

437    the phenotype is expressed or not.

438          Pangenome reconstruction for *Bo* and *Bx* revealed extensive variability between strains of

439    these closely related species, which is not unexpected for bacteria that engage in LGT. However,

440    the lack of mobile DNA signatures for the majority of accessory genes and evidence of inter-

441    genomic recombination between species at core genes provides new insight into what may be a

442    prominent mechanism of genome diversification in members of this phylum. The previously

443    described intergenomic transfer mechanisms in *B. thetaiotaomicron* and *B. fragilis* required the

444    presence of active or inactive ICEs, highlighting the potential roles for these mobile elements in

445    not just shaping genomes directly but also indirectly through their ability to catalyze exchange of

446    broader genomic segments. In *B. thetaiotaomicron*, genome transfer was determined to initiate at

447    genomically-integrated ICEs of which there are four in the type strain of *B. thetaiotaomicron*

448    (VPI-5482). These have not been shown to be fully functional for circularization and

449    mobilization. However, introduction and activation of an additional, excision-proficient

450    conjugative transposon (either cTnDOT or cTnERL) (42), which share common features with the

451    genomic ICEs, catalyzed expression of genes in the genomic ICEs and transfer of parts of the

452    genome in a manner that requires *recA* and homologous DNA to be present in the recipient (42).

453    An additional study in *B. fragilis* showed that conjugation from a strain with multiple genomic

454    ICEs, with one or more presumably retaining transfer activity, results in transfer of up to 435Kb

455    of chromosome into a recipient that initiates near genomic ICEs, with individual transfer events

456    being of variable size. The latter observation suggests that intergenomic recombination could

457    then occurs at different homologous regions (*i.e.*, the core gene nodes observed in the

458    pangenome), which could depend on the amount of genomic DNA transferred and the

459    length/homology of available recombination sites. Given that the number of ICEs in individual

460    genomes is variable, and their ability to be activated by functional conjugative transposons that

461    are circulating in the ecosystem may also vary, it will be interesting to determine in future work

462    if there are hotspots for genome transfer or if certain strains/species are dominant genome donors

463    that could play a disproportionate role.

464          The phenotypic similarity between members of the same species (*e.g.*, *Bo* and *Bx*) and the

465    large amount of gene diversity, including genes involved in carbohydrate metabolism, presents a

466    paradox and raises the question of why the genome diversification observed in strains of *Bo* and

467    *Bx* has not pushed members of these species to behave more differently and cluster based on

468    phenotype with members of other species. One answer may be the apparent exclusion of some

469    traits, such as mucin *O*-glycan/hemicellulose metabolism, which may limit the fitness advantage

470    associated with acquiring new phenotypes. A second emerges from the proposed genome-

471    exchange mechanism for which we offer new experimental support. Since this intergenomic

472    exchange relies on homologous recombination, its frequency should decrease between genomes

473    that are more divergent. Thus, this strategy may be one mechanism through which only closely

474    related bacteria can share traits that are advantageous with other close relatives. The presence of

475    orthologous PULs that confer the same function (*e.g.*, GluM and GalM utilization), some of

476     which appear to be subjected to LGT, further complicates interpretations of genotype-to-

477     phenotype relationships in these bacteria. Notably, the genome transfer mechanism proposed

478     here does not account for how new genes can be incorporated between conserved nodes. Rather,

479     this variability must pre-exist among different strains and therefore be created by different inter-

480     and intragenomic diversification mechanisms. Nevertheless, the data that we report here

481     underscore the notion that individual gut symbiont genomes are not just highly variable, but also

482     dynamically so.

483

484     **Materials and methods**

485     **Bacterial strains and growth conditions**

486     A total of 354 human and animal gut Bacteroidetes were included in this study. A

487     complete list is provided in **Table S1**, along with species designation based on 16S rRNA gene

488     sequencing and associated meta-data. Dr. Abigail Salyers (University of Illinois, Urbana-

489     Champagne) kindly provided many of the strains and two large portions of this collection were

490     isolated over several decades: 99 strains with "WH" designations were collected from fecal

491     samples of healthy human volunteers as part of the Woods Hole Summer Course on Microbial

492     Diversity in the late 1990s; 95 additional strains with "VPI" designations were collected from

493     human samples at the Virginia Polytechnic Institute in the 1960s-1970s. Species classifications

494     were made based on alignment of a minimum of 734 bp of 16S rRNA gene sequence to a

495     database containing the type strains of >29 named human gut Bacteroidetes species using the

496     classify.seqs command with Bayesian settings in the program mothur (46); assignment for each

497     strain was also manually checked by Blast (47). Isolates with ≥98% 16 rDNA gene sequence

498     identity to the type strain of a named species were labeled with that species designation. This

499    classification strategy included all except for three of the 354 strains examined, which ranged

500    between 96.6 to 96.7% sequence identity to the *B. uniformis* ATCC type strains and based on

501    sequential isolate numbers might be clones from the same individual (see WH15, WH16, WH17

502    entries in **Table S1**). Because of the small number of strains that did not satisfy our 98% cutoff,

503    we grouped these unclassified strains with their nearest relative and label them as more divergent

504    in **Table S1**; although, in most cases the carbohydrate phenotypes of these strains were very

505    similar to other members of the *B. uniformis* group.

506        All strains were routinely grown in an anaerobic chamber (Coy Lab Products, Grass

507    Lake, MI) at 37°C under an atmosphere of 5% $H_2$, 5% $CO_2$, and 90% $N_2$ on brain-heart infusion

508    (BHI, Beckton Dickinson) agar that included 10% defibrinated horse blood (Colorado Serum

509    Co.) and gentamicin (200 µg/ml). A single colony was picked into either tryptone-yeast extract-

510    glucose (TYG) media (48) or modified chopped-meat carbohydrate broth (**Table S8**) and then

511    sub-cultured into a minimal medium (MM) formulation that contained a mixture of

512    monosaccharides, vitamins, nucleotides, amino acids and trace minerals (**Table S2** provides

513    components and a complete recipe).

514

515    **Carbohydrate growth array setup and data collection**

516        Two different minimal medium formulations were used in the carbohydrate growth

517    arrays (**Table S1** lists the formulation used for each isolate). The simpler of the two formulations

518    (medium 1) was identical to the above MM, except that no carbohydrates were included and the

519    medium was prepared at 2X concentration. The second minimal medium formulation (medium

520    2) was identical to medium 1, but included beef extract (0.5% w/v final concentration) as an

521    additional supplement. We initially attempted to cultivate all of the species tested using only

522   medium 1, but determined that beef extract was specifically required to allow growth of some

523   species, especially *Parabacteroides* spp., *Barnesiella intestinihominis*, *Odoribacter splanchnicus*

524   and the branch of *Bacteroides* that includes *B. plebeius* and *B. massiliensis*. Growth in the

525   absence of an added carbohydrate source was generally not observed ot very low, except with

526   *Parabacteroides* that may be able to grow to a low level on the added 0.5% beef extract. The

527   corresponding negative control wells for each strain assayed were averaged and this value

528   subtracted from the total growth calculation of the corresponding to strain on other carbohydrates

529   tested. Despite several attempts to supplement minimal media with different components or

530   employ more stringent anaerobic methods, we were unable to cultivate several common

531   Bacteroidetes genera/species (*Prevotella* spp., *Paraprevotella* spp., *Alistipes* spp., and

532   *Bacteroides coprocola* and *Bacteroides coprophilus*) in these two MM formulations and

533   therefore did not include them in this study. All of these isolates readily grew in rich medium,

534   suggesting that they have specific nutritional requirements that were not met in the MM

535   formulations used.

536       Carbohydrate growth arrays were run as described previously (23) using a list of 45

537   carbohydrates (**Table S9)** that were present in duplicate, non-adjacent wells of a 96-well plate;

538   two additional wells contained no carbohydrate and served as negative controls. Each MM was

539   prepared as a 2X concentrated stock without carbohydrates (MM-no carb). An aliquot of each

540   strain was taken from a MM-monosaccharides culture (grown for 16-20 h) and was centrifuged

541   to pellet cells. Bacteria were resuspended in the same volume of 2X MM-no carb and then

542   centrifuged again prior to suspension in a volume of 2X MM-no carb that was equal to the

543   original volume. These washed bacterial cells were then inoculated at a 1:50 ratio into 2X MM-

544   no carb and the suspension was added in equal volume (100µl/well) to the 96 wells of the

545    carbohydrate growth array. Each well of the carbohydrate growth array contained 100µl of 2X

546    carbohydrate stock (10-20mg/ml); thus, when diluted 2-fold resulted in 1X MM containing a

547    unique carbohydrate and a bacterial inoculum that was identical to other wells. Growth arrays

548    were monitored at kinetic intervals of 10-20 minutes using a microplate stacking device and

549    coupled absorbance reader (Biotek Instruments; Winooski, VT) and data recorded for 4 d

550    (variable kinetic interval times reflect variations in the number of microtiter plates present in a

551    given batch).

552

553    **Carbohydrate growth array data processing**

554    Growth data were processed according to the following workflow: 1. data for each strain

555    were exported from Gen5 software (Biotek Instruments; Winooski, VT) into Microsoft Excel

556    and a previously described automated script was employed to call the points at which growth

557    began (min) and ended (max) (23); 2. Each file was manually checked to validate that

558    appropriate calls were made and the min and max values edited if needed (generally, only due to

559    obvious baselining artifacts or erroneously high calls caused by temporary bubbles or

560    precipitation); 3. "total growth" ($A_{600}$ max – $A_{600}$ min) and "growth rate" [($A_{600}$max – $A_{600}$ min) /

561    ($t$ max – $t$ min) were calculated for each strain on each substrate ($A_{600}$ is the absorbance value at

562    600 nm that corresponds to each min and max point; $t$ is the corresponding time values in

563    minutes); 4. Individual cultures in which total growth was ≤ 0.1 were scored as "no growth" and

564    their $A_{600}$ values converted to 0. Only assays in which both replicates showed an increase in $A_{600}$

565    ≥ 0.1 were considered as growth; if the two replicate assays were discordant (one positive, one

566    negative), then both values were converted to zero.

567       To normalize the results for each strain, the substrate(s) that provided maximum total

568    growth and growth rate values were determined and these were set to 1.0. All other growth

569    values for a given strain were normalized to this maximum value, providing a range of values

570    between 0 and 1.0. We next normalized growth ability across individual substrates using the

571    previously normalized values for each individual strain: the strain with the maximum total

572    growth and growth rate values were identified (many of these were already set to 1.0). Then, the

573    corresponding values for each other species on that particular substrate were calculated as a

574    fraction of the maximum value for that substrate, yielding a range of values between 0 and 1.0

575    for each substrate. These values were used to create the heat map shown in **Figs. 2** and **S3** and all

576    raw and normalized values are provided in **Table S1**.

577

578    **Data clustering and statistics**

579       Heatmaps and corresponding dendrograms were generated using the "heatmap" function

580    in the "stats" package of R (version 3.4.0) which employs unsupervised hierarchical clustering

581    (complete linkage method) to group similar carbohydrate growth profiles. Pearson Correlation

582    was used to calculate co-occurrence of the ability to grow on each pair of different substrates.

583    The normalized growth value for each substrate was compared to the corresponding growth

584    values on all other substrates using the Pearson correlation test in R and these values are

585    displayed in the Pearson correlation plot in **Figure S5**.

586

587    **Pangenome reconstruction for *B. ovatus* and *B. xylanisolvens* strains**

588       Since one of the seven strains used for pangenome reconstruction (*B. xylanisolvens*

589    XB1A) was assembled into a single circular chromosome, we used this genome as a scaffold for

590    the contigs representing the remaining six strains. Contigs from the six unfinished strains were

591    aligned against the XB1A genome using a combination of Mauve (50), to align and orient larger

592    contigs, and reciprocal best Blast-hit analysis using ≥90% amino acid identity to identify likely

593    homologs, to provide finer resolution. Contigs from draft genome assemblies or *Bx* XB1A were

594    broken as needed to accommodate the inclusion of unique accessory genes, but only in

595    circumstances where genes on both sides of the break could be aligned to homologs in one or

596    more genomes with a contig that spanned that break point. After constructing a preliminary

597    assembly, we analyzed the size distribution of putative homologous ORFs as a measure of

598    assembly accuracy and to identify variations in genetic organization that might be attributable to

599    real genetic differences such as frame shifts, which would result in two homologous gene calls of

600    smaller size in the genome containing the frameshift. Any variation >50% of homologous ORF

601    size was inspected manually using the "orthologous neighborhood viewer, by best Blast hit"

602    function in the U. S. Dept. of Energy Integrated Microbial Genomes (IMG) website. Introduced

603    contig breaks are documented in **Table S2** and **Fig. S6**. GenVision software (DNAstar, Madison,

604    WI) was used to visualize and label selected functions in the pan-genome assembly and also

605    display RNAseq data as a function of shared and unique PULs.

606

607    **RNAseq analysis**

608         For RNAseq, *B. xylanisolvens* and *B. ovatus* cells were grown to mid-exponential phase

609    on either purified mucin *O*-linked glycans (purified in house from Sigma Type III porcine gastric

610    mucin) or glucose as a reference as previously described (22). Total RNA was extracted using an

611    RNeasy kit (Qiagen), treated with Turbo DNase I (Ambion), and mRNA was enriched using the

612    Bacterial Ribo-Zero rRNA removal kit (Epicentre). Residual mRNA was converted to

613     sequencing libraries using TruSeq barcoded adaptors (Illumina) and sequenced at the University

614     of Michigan Sequencing Core in an Illumina HiSeq instrument with 24 samples multiplexed per

615     lane. Bar-coded data were demultiplexed and analyzed using the Arraystar software package

616     with Qseq (DNAstar). All RNAseq data are publicly available from the National Institutes of

617     Health Gene Expression Omnibus Database under accession numbers GSM4714867-

618     GSM4714890.

619

620     **Core gene determination and detection of LGT events between *Bo* and *Bx* strains**

621         The core gene alignment was generated with cognac (51). The alignment was then

622     partitioned into the individual component genes and approximate maximum likelihood gene trees

623     were generated with fastTree (52). Co-phylogenetic distances were calculated with APE (53). A

624     distance threshold of greater than 0.1 to the same species and less than 0.1 to the opposite species

625     was used to identify alleles bearing signatures of HGT. All analyses were performed in R

626     (version 3.6.3) (54). All code developed for this project are available at

627     https://github.com/rdcrawford/bacteroides_hgt.

628     **Acknowledgements**

640    **Competing Interests:** The authors declare no competing interests.

641

642    **Figure legends**

643    **Figure 1. Glycan degradation abilities among gut Bacteroidetes. (A.)** The number of species

644    out of 29 tested that degrade each polysaccharide is listed in order of decreasing degradation

645    frequency from left to right. Since not all strains within a given species necessarily have the

646    metabolic potential to utilize each polysaccharide, colors illustrate the percentage of strains

647    within each degrading species that possess the indicated ability.  **(B.)** The number of

648    polysaccharides that a given species degrades in decreasing order. The number of strains tested

649    for each species is listed in parentheses and colors represent the percentage of strains in each

650    indicated species that degrade each glycan counted towards the total.

651

652    **Figure 2. Heat map of individual polysaccharide utilization traits.** Species are clustered by

653    glycan utilization phenotype based on normalized total growth level (**Fig. S4B**). The magnitude

654    of growth is indicated by the heatmap scale at the bottom right. Columns at the left indicate the

655    source (human, animal) and time period of isolation. The cladogram at the far left shows the

656    results of unsupervised clustering of the data based the normalized growth data shown. The

657    species designations at the right are the results of 16S rRNA gene sequencing (>98% identity to

658    the species type strain was used to assign species). All raw and normalized growth and rate data

659    for individual strains may be found in **Table S1** see **Fig. S3** for an expanded heatmap with

660    monosaccharide data and individual strain names labels.

661

662    **Figure 3. Host mucin *O*-glycan metabolism within the *Bacteroides*. (A.)** A phylogenetic tree

663    based on housekeeping genes that compares mucin *O*-glycan utilization across species. The

664    diameter of the black circles represents the number of strains tested within each species (sample

665    depth), whereas the size of the overlaid red circle corresponds to the number of strains exhibiting

666    *O*-glycan metabolism. Note that some species have either full or no penetrance of this

667    phenotypic trait yet others like *B. ovatus*/*B. xylanisolvens* have more extensive variability among

668    strains. **(B.)** Strains of *B. ovatus* (blue) and *B. xylanisolvens* (green) that show variable growth

669    abilities on mucin *O*-glycan (n=2 growth assays per bar, error bars are range between values).

670    Gray histogram bars are total growth controls on an aggregate of the monosaccharides that all

671    strains of these two species grow on (**Table S1**) and are provided as a reference for overall

672    growth ability on a non *O*-glycan substrate. Data from two established *O*-glycan degraders, *B.*

673    *massiliensis* and *B. thetaiotaomicron,* are also shown for reference. Species with black arrows

674    were used for pangenome analyses to compare genetic traits associated with mucin *O*-glycan

675    metabolism. We performed RNA-seq on three strains included in this pangenome analyses

676    (black boxes) positive for *O*-glycan utilization and an additional strain, *B. ovatus* NLAE-zl-H59

677    (red arrow, box), to see if there were unique genes/PULs present in strains that have the ability to

678    grow on mucin *O*-glycans.

679

680     **Figure 4. Distribution of all genes as well as core polysaccharide utilization functions in the**

681     ***Bo/Bx* pangenome. (A.)** The left panel shows the number of core genes (*i.e*., those present in all

682     seven strains used for pangenome construction) compared to genes present in 2-7 of the

683     individual strains. The right panel shows the same distribution of genes assigned to PULs or

684     particular degradative CAZyme families (GH, PL, CE, see **Tables S2** and **S3** for more detailed

685     assignments. **(B.)** The distribution of genes between mucin-degrading (n=3) and non-degrading

686     (n=4) strains used to construct the pangenome. Top numbers indicate total genes, while numbers

687     in parentheses indicate the number of PULs (not individual PUL genes) in each category. **(C.)**

688     Distribution of the genes that are unique to the three mucin-degrading strains within each

689     genome. Genes/PULs are numbered as described for B. Note that no PULs are shared by all three

690     strains.

691     **Figure 5. (A.)** A higher-resolution view of a region of the *Bo/Bx* pangenome shows the variable

692     presence of at least six different PULs occurring between three genomic nodes (nodes 33-35 in

693     this quarter of the total pangenome. Segment 2 of the physical pangenome map was selected

694     because the first segment initiated with numerous small contigs and this segment contained

695     previously validated genes for xyloglucan metabolism (49). Node genes are colored red, while

696     *susC*-like and *susD*-like genes are colored purple and orange, respectively, and glycoside

697     hydrolase genes in light blue. GH family numbers are given below select PULs starting from the

698     top to indicate potential specificity and new numbers are only added going down the schematic if

699     the family assignments are different, indicating a different PUL. A well-studied *B. ovatus* PUL

700     for xyloglucan degradation (49) is shown in the center and occurs variably between two nodes

701     and also has variable gene content. The two bottom genomes are from different species,

702     *Bacteroides finegoldii* (*Bfin*) and *Bacteroides fragilis* (*Bfra*) and show less complex genome

703 architecture with the *Bfra* region possessing no PULs. **(B.)** A broader view of the genome region

704 shown in A. showing that the same mosaic pattern is common across the pangenome. Only PULs

705 are illustrated, although many other genes were also variable in these regions. The numbers at

706 the bottom delineate the presence of 35 different core gene nodes (as in A. some nodes contain

707 multiple core genes) in this section of the genome and the presence of homologous or unique

708 PULs is illustrated according to the color code at right (see **Fig. S6** for high resolution physical

709 maps of the pan-genome with PUL annotations). Note that in some cases up to five different

710 PULs were located at one location **(C.)** A schematic showing the proposed mechanism of

711 genome exchange based on previous studies (42-44) and observations presented here. Genomic

712 ICEs that are either partially active (excision deficient, but capable of initiating DNA strand

713 breakage and conjugation) or activated *in trans* by the presence of an exogenous conjugative

714 transposon, initiate genome mobilization from a donor into a recipient. If sufficient homology

715 between node genes exists in the recipient, homologous recombination between two nodes can

716 replace a section of the recipient with a segment from the donor. Note that genomic regions are

717 shown as linear fragments for simplicity, but would be circular.

718 **Figure 6. (A.)** Schematic of the workflow to identify putative LGT core genes: align genes and

719 build corresponding trees for each core gene,determine the median substitution distances

720 distances for each allele of a core gene in a given strain to both species, and identify loci with

721 an identical conserved structure between isolates of opposite species. **(B.)** Plot of median

722 distances for all core genes identified in the 33 genomes analyzed. The boxes show the regions

723 containing genes for which the median distance was $\geq 0.1$ to the assigned species for a given

724 strain and $\leq 0.1$ for the opposite species to which a strain is assigned. These genes were

725 determined to be high-confidence examples of core/node genes that had been replaced by an

726    allele from the other species. **(C.)** A region of the *Bo/Bx* pangenome that contains a PUL

727    involved in galactomannan (GalM) and glucomannan (GluM) degradation. This PUL is present

728    in six strains of *Bx* and two strains of *Bo* and in the latter cases flanking node genes exhibit

729    signatures of being derived from LGT with a *Bx* donor (the yellow box highlights potential

730    recombination region). The columns at the left indicate the growth of each strain on GalM or

731    GluM. The ability to grow on GalM is fully correlated with the presence of one of two different

732    PULs, or both, that are transcriptionally activated during growth on this substrate (**Fig. S9**) (23).

733    Notably, some strains (red "+") are able to grow weakly on GluM but do not possess either of the

734    identified PULs, suggesting that additional, partially orthologous PULs exist that confer the

735    ability to use only GluM.

736

737    **Supplementary Figure and Table Legends**

738    **Figure S1.** Schematics of the polysaccharides used in this study with sugar composition and

739    linkages schematized according to the "Symbol nomenclature for glycans" standard format and

740    based on the symbol key provided at the right. Linkages are labeled as α or β and the number

741    provided represents the carbon position in the recipient sugar. The carbon in the donor sugar is

742    carbon-1 in all cases except N-acetyl neuraminic acid and is not shown. Note that pectic galactan

743    (potato and lupin), xylan (oat spelt and wheat arabinoxylan) and amylopectin (potato and maize)

744    can have variable structures based on plant source. Abbreviations for several polysaccharides are

745    provided in parentheses and used throughout the text and figures.

746 **Figure S2. Correlation of replicate growth and rate measurements.** Two replicate

747 measurements were made for each of the two parameters recorded, total growth (**A.**) and growth

748 rate (**B.**) for each species on each carbohydrate substrate. Data points are color-coded based on

749 whether the two replicates exhibited variation between 0-5% (black), 5-10% (blue), 10-20%

750 (green), >20% (orange) or growth in one assay and no growth in the other (red). (**C.**) A linear

751 function was fitted (with red points omitted) to calculate an $r^2$ value for the data set associated

752 with utilization of each individual substrate. Measurements on some substrates were more

753 variable than on others due, at least in part, to the tendency of these substrates to partially

754 precipitate or retrograde during growth, which yielded variable levels of background absorbance.

755 **Figure S3.** A heatmap identical to the one shown in **Fig. 2** main text, except that

756 monosaccharide growth data is included. Strain names are also noted at the far right (best viewed

757 in electronic PDF form with magnification) and animal strains are labeled in red font.

758 **Figure S4. (A.)** A scheme for evaluating which aspects of growth phenotype data are most

759 influential for clustering strains that belong to the same species using hypothetical *B. theta* data

760 as an illustrative example. A quantitative index was used in which the number of strains tested is

761 divided by the minimum number of branches needed to encompass all of the strains for that

762 species, with a perfect score being "1" (*e.g.*, eight *B. theta* strains divided by the minimum of

763 eight branches needed to encompass all strains in the top example). **(B.)** Actual clustering index

764 data for the raw and normalized growth and rate data gathered for 354 different Bacteroidetes

765 strains. M and P stand for "monosaccharide" and "polysaccharide" growth, respectively. One of

766 the two most optimal conditions, which incorporates normalized growth data on polysaccharides

767 only, was used to construct **Figs. 2** and **S3**.

768     **Figure S5.** A Pearson correlation plot to determine if individual growth abilities co-occur in the

769     same strains. Positive or negative correlations that are ≥0.40 are shown in the colors indicated.

770     **Figure S6.** High-resolution maps of the entire reconstructed pangenome. These maps are

771     provided in four separate parts due to their large size (labeled as **Fig. S6a, b, c, d**) and

772     correspond to the data table provided in **Table S2**. Note that a 5$^{th}$ file is provided with

773     information about the gene, locus and strand breaking legend data.

774     **Figure S7**. Circular pangenome and corresponding mucin *O*-glycan transcriptomics from *Bx*

775     D22, *Bo* 3_1_23 and *Bo* D2 .

776     **Figure S8.** Individual maps of high-confidence inter-genomic exchange events between Bo and

777     Bx strains.

778     **Figure S9. (A.)** Schematics of PUL-A and PUL-B associated with GalM and GlcM utilization.

779     In *Bo* ATCC8384, elimination of PUL-A eliminates both of these growth abilities. (**B.**).

780     Expression analysis by qPCR of two sentinel genes from PUL-B in *Bo* strain D2 that lacks PUL-

781     A but still exhibits robust growth on GalM.

782     **Table S1.** Strain designations, growth levels, growth rates, host species, isolation periods,

783     growth media, 16S rRNA similarities and, if applicable, public genome sequence references for

784     all Bacteroidetes strains used in this study.

785

786     **Table S2.** Data table of the reconstructed pangenome of seven *Bo/Bx* strains. Additional notes

787     are provided directly on the table.

788

789    **Table S3.** PULs that were delineated in the seven strain pangenome with annotations based on

790    whether they were unique to mucin non-degrading strains, unique to mucin-degraders or shared

791    between strains in both categories. Additional notes are provided directly on the table.

792

793    **Table S4.** Gene expression changes detected using whole-genome transcriptional profiling by

794    RNA-seq of *B. xylanisolvens* D22 grown on mucin *O*-glycan as a sole carbon source compared

795    to glucose reference. Additional notes are provided directly on the table.

796

797    **Table S5.** Gene expression changes detected using whole-genome transcriptional profiling by

798    RNA-seq of *B. ovatus* 3-1-23 grown on mucin *O*-glycan as a sole carbon source compared to

799    glucose reference. Additional notes are provided directly on the table.

800

801    **Table S6.** Gene expression changes detected using whole-genome transcriptional profiling by

802    RNA-seq  of *B. ovatus* D2 grown on mucin *O*-glycan as a sole carbon source compared to

803    glucose reference. Additional notes are provided directly on the table.

804

805    **Table S7.** Gene expression changes detected using whole-genome transcriptional profiling by

806    RNA-seq  of *B. ovatus* NLAE-zl-H59 grown on mucin *O*-glycan as a sole carbon source

807    compared to glucose reference. Additional notes are provided directly on the table.

808

809    **Table S8.** Liquid media recipes (sheet **A**) and components (sheet **B**) for growing the

810    Bacteroidetes used in this study.

811

812    **Table S9.** Mono- and polysaccharides used in the phenotypic growth arrays and corresponding

813    supplier or purification details.

814

**References**

1. Koropatkin NM, Cameron EA, Martens EC. How glycan metabolism shapes the human gut microbiota. *Nat Rev Microbiol*. 2012; 10: 323-35.

2. Flint HJ, Scott KP, Duncan SH, Louis P, Forano E. Microbial degradation of complex carbohydrates in the gut. *Gut microbes*. 2012; 3:289-306.

3. Porter NT, Martens EC. The Critical Roles of Polysaccharides in Gut Microbial Ecology and Physiology. *Annual review of microbiology*. 2017; 71: 349-69.

4. El Kaoutari A, Armougom F, Gordon JI, Raoult D, Henrissat B. The abundance and variety of carbohydrate-active enzymes in the human gut microbiota. *Nature reviews Microbiology*. 2013; 11: 497-504.

5. McNeil NI. The contribution of the large intestine to energy supplies in man. *The American journal of clinical nutrition*. 1984; 39: 338-42.

6. Desai MS, Seekatz AM, Koropatkin NM, Kamada N, Hickey CA, Wolter M, et al. A Dietary Fiber-Deprived Gut Microbiota Degrades the Colonic Mucus Barrier and Enhances Pathogen Susceptibility. *Cell*. 2016; 167: 1339-53.

7. Sonnenburg ED, Smits SA, Tikhonov M, Higginbottom SK, Wingreen NS, Sonnenburg JL. Diet-induced extinctions in the gut microbiota compound over generations. *Nature*. 2016; 529: 212-5.

8. Sonnenburg ED, Sonnenburg JL. Starving our microbial self: the deleterious consequences of a diet deficient in microbiota-accessible carbohydrates. *Cell metabolism*. 2014; 20: 779-86.

9. Wrzosek L, Miquel S, Noordine ML, Bouet S, Joncquel Chevalier-Curt M, Robert V, et al. *Bacteroides thetaiotaomicron* and *Faecalibacterium prausnitzii* influence the production of mucus glycans and the development of goblet cells in the colonic epithelium of a gnotobiotic model rodent. *BMC Biol.* 2013; 11: 61.

10. Ganapathy V, Thangaraju M, Prasad PD, Martin PM, Singh N. Transporters and receptors for short-chain fatty acids as the molecular link between colonic bacteria and the host. *Current opinion in pharmacology*. 2013; 13: 869-74.

11. Cook SI, Sellin JH. Review article: short chain fatty acids in health and disease. *Aliment Pharmacol Ther.* 1998; 12: 499-507.

12. Smith PM, Howitt MR, Panikov N, Michaud M, Gallini CA, Bohlooly YM, et al. The microbial metabolites, short-chain fatty acids, regulate colonic Treg cell homeostasis. *Science*. 2013; 341: 569-73.

859    13.  Kim M, Qie Y, Park J, Kim CH. Gut Microbial Metabolites Fuel Host Antibody Responses.
860    *Cell Host Microbe*. 2016; 20: 202-14.
861
862    14.  Xu J, Mahowald MA, Ley RE, Lozupone CA, Hamady M, Martens EC, et al. Evolution of
863    Symbiotic Bacteria in the Distal Human Intestine. *Plos Biol*. 2007; 5: e156.
864
865    15.  Cantarel BL, Coutinho PM, Rancurel C, Bernard T, Lombard V, Henrissat B. The
866    Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics.
867    *Nucleic acids research*. 2009; 37: D233-D8.
868
869    16.  McNulty NP, Wu M, Erickson AR, Pan C, Erickson BK, Martens EC, et al. Effects of diet
870    on resource utilization by a model human gut microbiota containing *Bacteroides*
871    *cellulosilyticus* WH2, a symbiont with an extensive glycobiome. *Plos Biol*. 2013; 11:
872    e1001637.
873
874    17.  Salyers AA, Vercellotti JR, West SE, Wilkins TD. Fermentation of mucin and plant
875    polysaccharides by strains of *Bacteroides* from the human colon. *Appl Environ Microbiol*.
876    1977; 33: 319-22.
877
878    18.  Salyers AA, West SE, Vercellotti JR, Wilkins TD. Fermentation of mucins and plant
879    polysaccharides by anaerobic bacteria from the human colon. *Appl Environ Microbiol*. 1977;
880    34: 529-33.
881
882    19.  Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, et al. A human gut
883    microbial gene catalogue established by metagenomic sequencing. *Nature*. 2010; 464: 59-
884    65.
885
886    20.  Ley RE, Hamady M, Lozupone C, Turnbaugh PJ, Ramey RR, Bircher JS, et al. Evolution of
887    mammals and their gut microbes. *Science*. 2008; 320: 1647-51.
888
889    21.  Eckburg PB, Bik EM, Bernstein CN, Purdom E, Dethlefsen L, Sargent M, et al. Diversity
890    of the human intestinal microbial flora. *Science*. 2005; 308: 1635-8.
891
892    22.  Martens EC, Chiang HC, Gordon JI. Mucosal Glycan Foraging Enhances Fitness and
893    Transmission of a Saccharolytic Human Gut Bacterial Symbiont. *Cell Host Microbe*. 2008; 4:
894    447-57.
895
896    23.  Martens EC, Lowe EC, Chiang H, Pudlo NA, Wu M, McNulty NP, et al. Recognition and
897    Degradation of Plant Cell Wall Polysaccharides by Two Human Gut Symbionts. *Plos Biol*.
898    2011; 9: e1001221.
899
900    24.  Bloom SM, Bijanki VN, Nava GM, Sun L, Malvin NP, Donermeyer DL, et al. Commensal
901    *Bacteroides* species induce colitis in host-genotype-specific fashion in a mouse model of
902    inflammatory bowel disease. *Cell Host Microbe*. 2011; 9: 390-403.

903 25. Hickey CA, Kuhn KA, Donermeyer DL, Porter NT, Jin C, Cameron EA, et al. Colitogenic
904 *Bacteroides thetaiotaomicron* Antigens Access Host Immune Cells in a Sulfatase-Dependent
905 Manner via Outer Membrane Vesicles. *Cell Host Microbe*. 2015; 17: 672-80.
906
907 26. Mazmanian SK, Round JL, Kasper DL. A microbial symbiosis factor prevents intestinal
908 inflammatory disease. *Nature*. 2008; 453: 620-5.
909
910 27. Lagkouvardos I, Pukall R, Abt B, Foesel BU, Meier-Kolthoff JP, Kumar N, et al. The
911 Mouse Intestinal Bacterial Collection (miBC) provides host-specific insight into cultured
912 diversity and functional potential of the gut microbiota. *Nat Microbiol*. 2016; 1: 16131.
913
914 28. Browne HP, Forster SC, Anonye BO, Kumar N, Neville BA, Stares MD, et al. Culturing of
915 'unculturable' human microbiota reveals novel taxa and extensive sporulation. *Nature*.
916 2016; 533: 543-6.
917
918 29. Lagier JC, Khelaifia S, Alou MT, Ndongo S, Dione N, Hugon P, et al. Culture of previously
919 uncultured members of the human gut microbiota by culturomics. *Nat Microbiol*. 2016; 1:
920 16203.
921
922 30. Tramontano M, Andrejev S, Pruteanu M, Klunemann M, Kuhn M, Galardini M, et al.
923 Nutritional preferences of human gut bacteria reveal their metabolic idiosyncrasies. *Nat
924 Microbiol*. 2018; 3: 514-22.
925
926 31. Hehemann JH, Kelly AG, Pudlo NA, Martens EC, Boraston AB. Bacteria of the human gut
927 microbiome catabolize red seaweed glycans with carbohydrate-active enzyme updates
928 from extrinsic microbes. *Proceedings of the National Academy of Sciences of the United
929 States of America*. 2012; 109: 19786-91.
930
931 32. Tamura K, Hemsworth GR, Dejean G, Rogers TE, Pudlo NA, Urs K, et al. Molecular
932 mechanism by which prominent human gut Bacteroidetes utilize mixed-linkage β-glucans,
933 major health-promoting cereal polysaccharides. Cell Reports. 2017; 21: 417-430.
934
935 33. Shipman JA, Berleman JE, Salyers AA. Characterization of four outer membrane
936 proteins involved in binding starch to the cell surface of *Bacteroides thetaiotaomicron*. *J
937 Bacteriol*. 2000; 182: 5365-72.
938
939 34. Déjean G, Tamura K, Cabrera A, Jain N, Pudlo N, Holm Viborg A, et al. Synergy between
940 cell-surface glycosidases and glycan-binding proteins dictates the utilization of specific
941 beta(1,3)-glucans by human gut *Bacteroides*. *mBio*. 2020; 11: e00095-20.
942
943 35. Pudlo NA, Urs K, Kumar SS, German JB, Mills DA, Martens EC. Symbiotic Human Gut
944 Bacteria with Variable Metabolic Priorities for Host Mucosal Glycans. mBio. 2015; 6:
945 e01282-15.
946

36.  Despres J, Forano E, Lepercq P, Comtet-Marre S, Jubelin G, Chambon C, et al. Xylan degradation by the human gut *Bacteroides xylanisolvens* XB1A(T) involves two distinct gene clusters that are linked at the transcriptional level. *BMC genomics*. 2016; 17: 326.

37.  Sonnenburg ED, Zheng H, Joglekar P, Higginbottom SK, Firbank SJ, Bolam DN, et al. Specificity of polysaccharide use in intestinal *Bacteroides* species determines diet-induced microbiota alterations. *Cell*. 2010; 141: 1241-52.

38.  Terrapon N, Lombard V, Drula E, Lapebie P, Al-Masaudi S, Gilbert HJ, et al. PULDB: the expanded database of Polysaccharide Utilization Loci. *Nucleic Acids Res*. 2018; 46:D677-D83.

39.  Johansson ME, Holmen Larsson JM, Hansson GC. Microbes and Health Sackler Colloquium: The two mucus layers of colon are organized by the MUC2 mucin, whereas the outer layer is a legislator of host-microbial interactions. *Proceedings of the National Academy of Sciences of the United States of America*. 2011; 108 Suppl 1: 4659-65.

40.  Pudlo NA, Vasconcelos-Pereira G, Parnami J, Cid M, Markert S, Tingley JP, et al. Extensive transfer of genes for edible seaweed digestion from marine to human gut bacteria. *Preprint*. 2021; BIORXIV/2020/142968.

41.  Martens EC, Kelly AG, Tauzin AS, Brumer H. The devil lies in the details: how variations in polysaccharide fine-structure impact the physiology and evolution of gut microbes. *Journal of molecular biology*. 2014; 426: 3851-65.

42.  Moon K, Sonnenburg J, Salyers AA. Unexpected effect of a *Bacteroides* conjugative transposon, CTnDOT, on chromosomal gene expression in its bacterial host. *Mol Microbiol*. 2007; 64: 1562-71.

43.  Husain F, Tang K, Veeranagouda Y, Boente R, Patrick S, Blakely G, et al. Novel large-scale chromosomal transfer in *Bacteroides fragilis* contributes to its pan-genome and rapid environmental adaptation. *Microb Genom*. 2017; 3:e000136.

44.  Whittle G, Hamburger N, Shoemaker NB, Salyers AA. A Bacteroides conjugative transposon, CTnERL, can transfer a portion of itself by conjugation without excising from the chromosome. *J Bacteriol*. 2006; 188: 1169-74.

45.  Reddy SK, Bagenholm V, Pudlo NA, Bouraoui H, Koropatkin NM, Martens EC, et al. A beta-mannan utilization locus in *Bacteroides ovatus* involves a GH36 alpha-galactosidase active on galactomannans. *FEBS Lett*. 2016; 590: 2106-18.

46.  Schloss PD, Westcott SL, Ryabin T, Hall JR, Hartmann M, Hollister EB, et al. Introducing mothur: open-source, platform-independent, community-supported software for describing and comparing microbial communities. *Appl Environ Microbiol*. 2009; 75: 7537-41.

992    47.  Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic Local Alignment Search
993    Tool. *J Mol Biol*. 1990; 215: 403-10.
994
995    48.  Holdeman LV, Cato ED, Moore WEC. Anaerobe Laboratory Manual. Moore WEC, editor.
996    Blacksburg, Va.: Virginia Polytechnic Institute and State University Anaerobe Laboratory;
997    1977.
998
999    49.  Larsbrink J, Rogers TE, Hemsworth GR, McKee LS, Tauzin AS, Spadiut O, et al. A
1000   discrete genetic locus confers xyloglucan metabolism in select human gut Bacteroidetes.
1001   *Nature*. 2014;506(7489):498-502.
1002
1003   50. Darling AC, Mau B, Blattner FR, Perna NT. Mauve: multiple alignment of conserved
1004   genomic sequence with rearrangements. *Genome Res*. 2004; 14: 1394-403.
1005
1006   51. Crawford RD, Snitkin ES. cognac: rapid generation of concatenated gene alignments for
1007   phylogenetic inference from large, bacterial whole genome sequencing datasets. *BMC*
1008   *Bioinformatics*. 2021; 22: 70.
1009
1010   52. Price MN, Dehal PS, Arkin AP. FastTree: computing large minimum evolution trees with
1011   profiles instead of a distance matrix. *Mol Biol Evol.* 2009; 26: 1641-50.
1012
1013   53. Popescu AA, Huber KT, Paradis E. ape 3.0: New tools for distance-based phylogenetics
1014   and evolutionary analysis in R. *Bioinformatics*. 2012; 28: 1536-7.
1015
1016   54. R Core Team. R: A language and environment for statistical computing. R Foundation
1017   for Statistical Computing, Vienna, Austria. 2017; URL https://www.R-project.org/.
1018

1019

1020

Figure 1

Figure 2

## Figure 3

Figure 4

## Figure 5

Figure 6

Figure S1

# Figure S2

**A.** Scatterplot of total growth replicates $R^2 = 0.95$

**B** Scatterplot of growth rate replicates $R^2 = 0.86$



Color codes:
● 0-5% variation
● 5-10% variation
● 10-20% variation
● >20% variation
● no growth in one replicate

**C.**

**Polysaccharides:**

| Substrate | $R^2$ between growth values | $R^2$ between rate values |
|---|---|---|
| AG | 0.99 | 0.89 |
| alg | 0.96 | 0.96 |
| α-mann | 0.93 | 0.76 |
| APm | 0.94 | 0.92 |
| APpo | 0.96 | 0.91 |
| arab | 0.98 | 0.96 |
| BBG | 0.95 | 0.70 |
| carr | 0.93 | 0.96 |
| Cell | 0.96 | 0.81 |
| CS | 0.96 | 0.91 |
| dex | 0.96 | 0.87 |
| GalM | 0.96 | 0.98 |
| GlcM | 0.93 | 0.85 |
| glyc | 0.96 | 0.94 |
| hep | 0.96 | 0.85 |
| hya | 0.91 | 0.88 |
| inulin | 0.92 | 0.89 |
| lam | 0.96 | 0.96 |
| levan | 0.96 | 0.88 |
| lich | 0.80 | 0.45 |
| MOG | 0.98 | 0.97 |
| OSX | 0.93 | 0.83 |
| PGA | 0.97 | 0.89 |
| PGI | 0.95 | 0.96 |
| PGp | 0.92 | 0.92 |
| por | 0.85 | 0.85 |
| pull | 0.84 | 0.78 |
| RGI | 0.96 | 0.98 |
| WAX | 0.97 | 0.42 |
| XyG | 0.92 | 0.72 |

**Monosaccharides:**

| Substrate | $R^2$ between growth values |
|---|---|
| Ara | 0.90 |
| Fru | 0.90 |
| Fuc | 0.93 |
| Gal | 0.60 |
| GalA | 0.86 |
| GalNAc | 0.93 |
| Glc | 0.69 |
| GlcA | 0.87 |
| GlcNAc | 0.72 |
| GlcNH3 | 0.93 |
| Man | 0.88 |
| NeuNAc | 0.86 |
| Rha | 0.96 |
| Rib | 0.94 |
| Xyl | 0.85 |

Figure S3

# Figure S4

**A.** Example of Cluster scoring scheme

**B.** Normalized data



Perfect grouping: 8 B. theta strains / 8 taxa on minimum branch = 1

Intermediate grouping: 8 B. theta strains / 16 taxa on minimum branch = 0.5

Weakest grouping: 8 B. theta strains / 21 taxa on minimum branch = 0.38

**C.** Unnormalized (raw) or binary (growth/no growth) data

P = polysaccharide data
M = monosaccharide data

|  | starches | | | | fructans | | GAGs | | | pectins | | | | | | hemi-celluloses | | | | | | | microbial & marine | | | | | | | O-glycans |
|  | Pull | Glyc | APp | APm | Inulin | Levan | Hep | Hya | CS | PGA | RGI | PGp | PGl | AG | Arab | GalM | GlcM | XyG | OSX | WAX | BBG | Cell | Lam | Lich | Dex | αmann | Alg | Carr | Porph | MOG |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pull | 1.00 | 0.60 | 0.43 | 0.51 | 0.10 | 0.11 | 0.14 | -0.03 | -0.06 | -0.04 | 0.10 | -0.07 | -0.12 | 0.03 | 0.07 | 0.05 | 0.04 | 0.04 | 0.07 | 0.02 | 0.10 | -0.12 | -0.14 | 0.02 | 0.12 | 0.08 | 0.09 | 0.02 | -0.12 | -0.05 |
| Glyc | 0.60 | 1.00 | 0.52 | 0.68 | 0.04 | 0.39 | 0.47 | 0.35 | 0.38 | 0.37 | 0.41 | 0.26 | 0.23 | 0.33 | 0.25 | 0.21 | 0.22 | 0.23 | 0.38 | 0.31 | 0.32 | 0.14 | 0.04 | 0.14 | 0.50 | 0.32 | 0.18 | 0.03 | -0.09 | -0.16 |
| App | 0.43 | 0.52 | 1.00 | 0.72 | 0.02 | 0.35 | 0.30 | 0.30 | 0.37 | 0.30 | 0.37 | 0.17 | 0.14 | 0.28 | 0.27 | 0.07 | 0.06 | 0.03 | 0.23 | 0.17 | 0.28 | 0.01 | -0.13 | 0.01 | 0.36 | 0.28 | 0.19 | 0.05 | -0.07 | -0.05 |
| Apm | 0.51 | 0.68 | 0.72 | 1.00 | -0.04 | 0.25 | 0.41 | 0.31 | 0.33 | 0.37 | 0.39 | 0.25 | 0.20 | 0.32 | 0.30 | 0.21 | 0.16 | 0.17 | 0.30 | 0.21 | 0.32 | 0.07 | -0.05 | 0.11 | 0.43 | 0.29 | 0.13 | 0.01 | -0.08 | -0.22 |
| Inulin | 0.10 | 0.04 | 0.02 | -0.04 | 1.00 | -0.03 | 0.09 | 0.16 | 0.08 | 0.08 | 0.00 | -0.07 | -0.01 | -0.02 | -0.23 | -0.01 | 0.00 | -0.02 | 0.06 | 0.09 | 0.08 | 0.10 | 0.00 | 0.08 | 0.05 | 0.03 | 0.09 | 0.03 | 0.07 | 0.05 |
| Levan | 0.11 | 0.39 | 0.35 | 0.25 | -0.03 | 1.00 | 0.62 | 0.55 | 0.63 | 0.64 | 0.58 | 0.44 | 0.43 | 0.53 | 0.37 | 0.08 | 0.09 | 0.14 | 0.31 | 0.29 | 0.20 | 0.09 | 0.12 | -0.04 | 0.55 | 0.47 | 0.33 | 0.10 | -0.04 | -0.07 |
| Hep | 0.14 | 0.47 | 0.30 | 0.41 | 0.09 | 0.62 | 1.00 | 0.65 | 0.69 | 0.66 | 0.66 | 0.38 | 0.36 | 0.35 | 0.26 | 0.18 | 0.17 | 0.16 | 0.54 | 0.52 | 0.47 | 0.35 | 0.04 | 0.00 | 0.70 | 0.37 | 0.34 | 0.06 | -0.04 | -0.22 |
| Hya | -0.03 | 0.35 | 0.30 | 0.31 | 0.16 | 0.55 | 0.65 | 1.00 | 0.83 | 0.77 | 0.69 | 0.44 | 0.45 | 0.44 | 0.26 | 0.05 | 0.02 | 0.04 | 0.40 | 0.38 | 0.40 | 0.23 | -0.09 | -0.08 | 0.64 | 0.56 | 0.37 | 0.07 | 0.05 | -0.10 |
| CS | -0.06 | 0.38 | 0.37 | 0.33 | 0.08 | 0.63 | 0.69 | 0.83 | 1.00 | 0.83 | 0.68 | 0.50 | 0.51 | 0.48 | 0.32 | 0.08 | 0.06 | 0.04 | 0.46 | 0.45 | 0.38 | 0.23 | -0.02 | -0.06 | 0.67 | 0.54 | 0.30 | 0.10 | 0.03 | -0.13 |
| PGA | -0.04 | 0.37 | 0.30 | 0.37 | 0.08 | 0.64 | 0.66 | 0.77 | 0.83 | 1.00 | 0.72 | 0.55 | 0.55 | 0.55 | 0.46 | 0.03 | -0.01 | 0.02 | 0.44 | 0.44 | 0.24 | 0.14 | -0.07 | -0.06 | 0.62 | 0.56 | 0.28 | 0.10 | 0.06 | -0.10 |
| RGI | 0.10 | 0.41 | 0.37 | 0.39 | 0.00 | 0.58 | 0.66 | 0.69 | 0.68 | 0.72 | 1.00 | 0.43 | 0.39 | 0.39 | 0.46 | 0.12 | 0.07 | 0.07 | 0.57 | 0.52 | 0.44 | 0.25 | -0.19 | -0.06 | 0.56 | 0.45 | 0.41 | 0.04 | 0.10 | -0.14 |
| PGp | -0.07 | 0.26 | 0.17 | 0.25 | -0.07 | 0.44 | 0.38 | 0.44 | 0.50 | 0.55 | 0.43 | 1.00 | 0.92 | 0.53 | 0.40 | 0.30 | 0.23 | 0.44 | 0.24 | 0.23 | 0.23 | 0.17 | 0.22 | 0.05 | 0.57 | 0.41 | 0.19 | 0.11 | -0.01 | -0.07 |
| PGl | -0.12 | 0.23 | 0.14 | 0.20 | -0.01 | 0.43 | 0.36 | 0.45 | 0.51 | 0.55 | 0.39 | 0.92 | 1.00 | 0.53 | 0.36 | 0.29 | 0.22 | 0.42 | 0.22 | 0.21 | 0.22 | 0.19 | 0.22 | 0.06 | 0.56 | 0.43 | 0.21 | 0.11 | -0.01 | -0.02 |
| AG | 0.03 | 0.33 | 0.28 | 0.32 | -0.02 | 0.53 | 0.35 | 0.44 | 0.48 | 0.55 | 0.39 | 0.53 | 0.53 | 1.00 | 0.54 | -0.03 | -0.07 | -0.02 | -0.05 | -0.08 | -0.06 | -0.13 | 0.15 | -0.03 | 0.48 | 0.66 | 0.04 | 0.05 | 0.02 | 0.22 |
| Arab | 0.07 | 0.25 | 0.27 | 0.30 | -0.23 | 0.37 | 0.26 | 0.26 | 0.32 | 0.46 | 0.46 | 0.40 | 0.36 | 0.54 | 1.00 | -0.03 | -0.14 | -0.04 | 0.05 | 0.02 | -0.22 | -0.25 | 0.04 | -0.02 | 0.24 | 0.45 | -0.02 | 0.03 | 0.05 | 0.20 |
| GalM | 0.05 | 0.21 | 0.07 | 0.21 | -0.01 | 0.08 | 0.18 | 0.05 | 0.08 | 0.03 | 0.12 | 0.30 | 0.29 | -0.03 | -0.03 | 1.00 | 0.81 | 0.74 | 0.42 | 0.41 | 0.52 | 0.51 | 0.30 | 0.24 | 0.38 | -0.13 | 0.18 | 0.05 | 0.16 | -0.42 |
| GlcM | 0.04 | 0.22 | 0.06 | 0.16 | 0.00 | 0.09 | 0.17 | 0.02 | 0.06 | -0.01 | 0.07 | 0.23 | 0.22 | -0.07 | -0.14 | 0.81 | 1.00 | 0.71 | 0.42 | 0.39 | 0.58 | 0.62 | 0.35 | 0.30 | 0.41 | -0.17 | 0.22 | 0.02 | -0.03 | -0.49 |
| XyG | 0.04 | 0.23 | 0.03 | 0.17 | -0.02 | 0.14 | 0.16 | 0.04 | 0.04 | 0.02 | 0.07 | 0.44 | 0.42 | -0.02 | -0.04 | 0.74 | 0.71 | 1.00 | 0.42 | 0.38 | 0.50 | 0.53 | 0.34 | 0.21 | 0.41 | -0.04 | 0.29 | 0.02 | -0.03 | -0.44 |
| OSX | 0.07 | 0.38 | 0.23 | 0.30 | 0.06 | 0.31 | 0.54 | 0.40 | 0.46 | 0.44 | 0.57 | 0.24 | 0.22 | -0.05 | 0.05 | 0.42 | 0.42 | 0.42 | 1.00 | 0.92 | 0.69 | 0.63 | -0.08 | 0.07 | 0.51 | 0.01 | 0.51 | 0.03 | 0.14 | -0.53 |
| WAX | 0.02 | 0.31 | 0.17 | 0.21 | 0.09 | 0.29 | 0.52 | 0.38 | 0.45 | 0.44 | 0.52 | 0.23 | 0.21 | -0.08 | 0.02 | 0.41 | 0.39 | 0.38 | 0.92 | 1.00 | 0.62 | 0.62 | -0.07 | 0.04 | 0.47 | -0.01 | 0.50 | -0.01 | 0.13 | -0.50 |
| BBG | 0.10 | 0.32 | 0.28 | 0.32 | 0.08 | 0.20 | 0.47 | 0.40 | 0.38 | 0.24 | 0.44 | 0.23 | 0.22 | -0.06 | -0.22 | 0.52 | 0.58 | 0.50 | 0.69 | 0.62 | 1.00 | 0.70 | -0.03 | 0.13 | 0.57 | -0.01 | 0.48 | 0.02 | -0.03 | -0.55 |
| Cell | -0.12 | 0.14 | 0.01 | 0.07 | 0.10 | 0.09 | 0.35 | 0.23 | 0.23 | 0.14 | 0.25 | 0.17 | 0.19 | -0.13 | -0.25 | 0.51 | 0.62 | 0.53 | 0.63 | 0.62 | 0.70 | 1.00 | 0.27 | 0.14 | 0.47 | -0.09 | 0.39 | -0.01 | -0.05 | -0.55 |
| Lam | -0.14 | 0.04 | -0.13 | -0.05 | 0.00 | 0.12 | 0.04 | -0.09 | -0.02 | -0.07 | -0.19 | 0.22 | 0.22 | 0.15 | 0.04 | 0.30 | 0.35 | 0.34 | -0.08 | -0.07 | -0.03 | 0.27 | 1.00 | 0.17 | 0.21 | 0.01 | -0.19 | 0.07 | 0.13 | -0.20 |
| Lich | 0.02 | 0.14 | 0.01 | 0.11 | 0.08 | -0.04 | 0.00 | -0.08 | -0.06 | -0.06 | -0.06 | 0.05 | 0.06 | -0.03 | -0.02 | 0.24 | 0.30 | 0.21 | 0.07 | 0.04 | 0.13 | 0.14 | 0.17 | 1.00 | 0.10 | -0.07 | -0.05 | -0.01 | -0.01 | -0.16 |
| Dex | 0.12 | 0.50 | 0.36 | 0.43 | 0.05 | 0.55 | 0.70 | 0.64 | 0.67 | 0.62 | 0.56 | 0.57 | 0.56 | 0.48 | 0.24 | 0.38 | 0.41 | 0.41 | 0.51 | 0.47 | 0.57 | 0.47 | 0.21 | 0.10 | 1.00 | 0.42 | 0.32 | 0.08 | -0.06 | -0.36 |
| αmann | 0.08 | 0.32 | 0.28 | 0.29 | 0.03 | 0.47 | 0.37 | 0.56 | 0.54 | 0.56 | 0.45 | 0.41 | 0.43 | 0.66 | 0.45 | -0.13 | -0.17 | -0.04 | 0.01 | -0.01 | -0.01 | -0.09 | 0.01 | -0.07 | 0.42 | 1.00 | 0.11 | 0.13 | -0.03 | 0.19 |
| Alg | 0.09 | 0.18 | 0.19 | 0.13 | 0.09 | 0.33 | 0.34 | 0.37 | 0.30 | 0.28 | 0.41 | 0.19 | 0.21 | 0.04 | -0.02 | 0.18 | 0.22 | 0.29 | 0.51 | 0.50 | 0.48 | 0.39 | -0.19 | -0.05 | 0.32 | 0.11 | 1.00 | -0.03 | -0.02 | -0.24 |
| Carr | 0.02 | 0.03 | 0.05 | 0.01 | 0.03 | 0.10 | 0.06 | 0.07 | 0.10 | 0.10 | 0.04 | 0.11 | 0.11 | 0.05 | 0.03 | 0.05 | 0.02 | 0.02 | 0.03 | -0.01 | 0.02 | -0.01 | 0.07 | -0.01 | 0.08 | 0.13 | -0.03 | 1.00 | 0.00 | -0.01 |
| Porph | -0.12 | -0.09 | -0.07 | -0.08 | 0.07 | -0.04 | -0.04 | 0.05 | 0.03 | 0.06 | 0.10 | -0.01 | -0.01 | 0.02 | 0.05 | 0.16 | -0.03 | -0.03 | 0.14 | 0.13 | -0.03 | -0.05 | 0.13 | -0.01 | -0.06 | -0.03 | -0.02 | 0.00 | 1.00 | 0.00 |
| MOG | -0.05 | -0.16 | -0.05 | -0.22 | 0.05 | -0.07 | -0.22 | -0.10 | -0.13 | -0.10 | -0.14 | -0.07 | -0.02 | 0.22 | 0.20 | -0.42 | -0.49 | -0.44 | -0.53 | -0.50 | -0.55 | -0.55 | -0.20 | -0.16 | -0.36 | 0.19 | -0.24 | -0.01 | 0.00 | 1.00 |

Color key:

- **perfect correlation (1.0)**
- **+0.7 to 1.0**
- **+0.40 to 0.7**
- -0.40 to 0.40
- **less than -0.40**

Figure S6

Please note that Figure S6 is provided as a zipped folder containing 4 separate quarters of the pangenome assembly along with a corresponding legend that explains the color coding scheme. The four maps correspond to **Table S2**.

Each map contains 5 vertically stacked panes of pangenome map starting in the upper left. Each horizontal pane has 8 rows with the top row representing the pangenome and the corresponding 7 individual genome regions shown below.

The example below shows a small region of pangenome section 2 in which *Bo* 3-1-23 is missing an ECF-σ regulated PUL. The small text above genes in individual genomes correspond to the contig and the green dashed line represents a region that was broken to accomodate accessory genes.

Figure S7



Legend:
- Transcriptome
- Genome/Pangenome
- Non PUL genes
- Unique PUL genes in mucin degraders
- Shared PUL genes between at least one mucin degrader and one non-degrader

# Figure S8. Bacteroides LGT Loci



HGT gene

Present

Absent

Allele more like B. ovatus

Allele more like B. xylanisolvens

# PUL LGT Events

## B. ovatus PUL LGT Event 1



species

TonB dependent receptor
fec operon regulator FecR
ECF RNA polymerase sigma factor SigK
Lipoprotein-releasing system transmembrane protein LolE
Ribosome-binding factor A
Putative O-methyltransferase/MSMEI_4947
Pyruvate kinase
3-dehydroquinate dehydratase
Tyrosine recombinase XerD
Tetratricopeptide repeat protein
Thiol-disulfide oxidoreductase ResA
Competence protein ComM
Aerobic respiration control sensor protein ArcB
hypothetical protein
SusD family protein
TonB dependent receptor
IPT/TIG domain protein
Alpha-xylosidase
Extracellular exo-alpha-(1->5)-L-arabinofuranosidase ArbA precursor
hypothetical protein
hypothetical protein
SusD family protein
TonB dependent receptor
hypothetical protein
hypothetical protein
Kojibiose phosphorylase
hypothetical protein
Sensor histidine kinase TodS

species
Bx
Bo

**B. ovatus PUL LGT Event 2**

**B. ovatus PUL LGT Event 3**

**B. ovatus PUL LGT Event 4**

**B. xylanisolvens PUL LGT Event 1**

## B. xylanisolvens PUL LGT Event 2

## B. xylanisolvens PUL LGT Event 3



species

Vitamin B12 dependent methionine synthase%2C activation domain
Uroporphyrinogen decarboxylase
Ribose-5-phosphate isomerase B
Transketolase
Intracellular exo-alpha-L-arabinofuranosidase 2
hypothetical protein
Xylulose kinase
L-arabinose isomerase
L-ribulose-5-phosphate 4-epimerase UlaF
bifunctional nicotinamide mononucleotide adenylyltransferase/ADP-ribose pyrophosphatase
Sodium/glucose cotransporter
Aldose 1-epimerase precursor
Extracellular exo-alpha-L-arabinofuranosidase precursor
Galactokinase
L-fucose-proton symporter
Aldose 1-epimerase precursor
putative mannose-6-phosphate isomerase GmuF
Tyrosine recombinase XerD
transcriptional activator RfaH
Polysaccharide biosynthesis/export protein
N-acetylmuramoyl-L-alanine amidase
hypothetical protein
hypothetical protein
hypothetical protein
hypothetical protein
hypothetical protein VirE N-terminal domain
hypothetical protein
UDP-glucose:undecaprenyl-phosphate glucose-1-phosphate transferase
Polysaccharide biosynthesis/export protein
Tyrosine-protein kinase ptk
putative glycosyltransferase EpsJ
hypothetical protein
putative acyl transferase

species
Bx
Bo

# non-PUL LGT Events

**B. ovatus non-PUL LGT Event 1**



species

species
Bx
Bo

Exoenzyme S synthesis regulatory protein ExsA
TonB dependent receptor
SusD family protein
Glycosyl hydrolase family 92
Glycosyl hydrolase family 92
Plant Basic Secretory Protein
Glycosyl hydrolase family 92
Retaining alpha-galactosidase precursor
Sensory transduction protein LytR
Sensor histidine kinase YehU
multidrug resistance protein MdtN
Cobalt-zinc-cadmium resistance protein CzcA
hypothetical protein Glutathione peroxidase homolog BsaA
Isopentenyl-diphosphate Delta-isomerase
S-adenosylmethionine:tRNA ribosyltransferase-isomerase
Arabinose operon regulatory protein
3-dehydroquinate synthase
hypothetical protein
hypothetical protein
Major cardiolipin synthase ClsA
Ribosomal RNA small subunit methyltransferase D
hypothetical protein
hypothetical protein
ATP-dependent RecD-like DNA helicase
hypothetical protein Bacteroidetes-Associated Carbohydrate-binding Often N-terminal
Endo-beta-N-acetylglucosaminidase F1 precursor
hypothetical protein
hypothetical protein
Susd and RagB outer membrane lipoprotein
TonB-dependent Receptor Plug Domain protein
hypothetical protein
hypothetical protein
Endo-beta-N-acetylglucosaminidase F1 precursor
F5/8 type C domain protein
Glycosyl hydrolase family 92
Glycosyl hydrolase family 92
fec operon regulator FecR
ECF RNA polymerase sigma factor SigW
Glycosyl hydrolase family 92
Alanine--tRNA ligase
Murein DD-endopeptidase MepM
HTH-type transcriptional repressor YcgE
Guanosine-3'%2C5'-bis(diphosphate) 3'-pyrophosphohydrolase
Membrane-bound lytic murein transglycosylase D precursor
hypothetical protein
putative chromosome-partitioning protein ParB
Sporulation initiation inhibitor protein Soj
hypothetical protein
5'-nucleotidase SurE
Lipid-A-disaccharide synthase
NigD-like protein
Phosphatidate cytidylyltransferase
ATP-dependent zinc metalloprotease FtsH
Ribosomal silencing factor RsfS
hypothetical protein
Magnesium transporter MgtE
Ribosomal RNA small subunit methyltransferase A
hypothetical protein
Cytosol non-specific dipeptidase

Bxylanisolvens_NLAE-zl-P727
Bxylanisolvens_NLAE-zl-P736
Bxylanisolvens_NLAE-zl-P593
Bxylanisolvens_NLAE-zl-P752
Bxylanisolvens_NLAE-zl-P792
Bxylanisolvens_CL03T12C04
Bacteroides_1_1_30
Bacteroides_2_1_22
Bacteroides_D1
Bacteroides_ovatus_SD_CC_2a
Bacteroides_xylanisolvens_SD_CC_1b
Bacteroides_xylanisolvens_XB1A
Bacteroides_D22
Bxylanisolvens_NLAE-zl-C29
Bxylanisolvens_NLAE-zl-H194
Bxylanisolvens_NLAE-zl-G310
Bxylanisolvens_NLAE-zl-G421
Bxylanisolvens_NLAE-zl-C182
Bxylanisolvens_NLAE-zl-C339
Bacteroides_2_2_4
Bacteroides_D2
Bovatus_NLAE-zl-C11
Bovatus_NLAE-zl-C34
Bacteroides_ovatus_CL03T12C18
Bovatus_NLAE-zl-H59
Bovatus_NLAE-zl-H73
Bacteroides_ovatus_ATCC_8483
Bovatus_C1D2T12C04
Bacteroides_3_8_47FAA
Bovatus_ovatus_3_1_33
Bacteroides_3_1_33
Bovatus_NLAE-zl-SD_CMC_3f
Bovatus_NLAE-zl-H304
Bovatus_NLAE-zl-H361

9

**B. ovatus non-PUL LGT Event 2**



species
Bx
Bo

**B. ovatus non-PUL LGT Event 3**

**B. ovatus non-PUL LGT Event 4**

**B. ovatus non-PUL LGT Event 5**



13

**B. ovatus non-PUL LGT Event 6**

**B. ovatus non-PUL LGT Event 7**

**B. ovatus non-PUL LGT Event 8**

**B. ovatus non-PUL LGT Event 9**

**B. xylanisolvens non-PUL LGT Event 1**



species

| | |
|---|---|
| | Bx |
| | Bo |

Glycosyl hydrolases family 2%2C sugar binding domain
hypothetical protein
hypothetical protein
Beta-galactosidase
Sensor histidine kinase TmoS
hypothetical protein
hypothetical protein
RNA polymerase sigma factor SigM
hypothetical protein
Inner membrane protein YbjJ
Inner membrane protein YjjP
hypothetical protein
hypothetical protein
3-oxoacyl-[acyl-carrier-protein] reductase FabG
Alpha-D-glucose-1-phosphate phosphatase YihX
hypothetical protein
hypothetical protein
hypothetical protein
N(omega)-hydroxy-L-arginine amidinohydrolase
putative HTH-type transcriptional regulator YybR
RNA polymerase sigma factor CarQ
fec operon regulator FecR
hypothetical protein Secretin and TonB N terminus short domain
metal-dependent hydrolase
Putative electron transport protein YccM
putative oxidoreductase
hypothetical protein
Vitamin B12 transporter BtuB precursor
hypothetical protein
Sensor histidine kinase TmoS
ABC-2 family transporter protein
Inner membrane transport permease YhhJ
Multidrug export protein EmrA
Outer membrane efflux protein
hypothetical protein
Proline--tRNA ligase
Alkaline phosphatase synthesis transcriptional regulatory protein PhoP
Sensor protein QseC
ABC-2 family transporter protein
putative ABC transporter ATP-binding protein YxlF
NPCBM-associated%2C NEW3 domain of alpha-galactosidase
hypothetical protein
hypothetical protein
hypothetical protein
tRNA N6-adenosine threonylcarbamoyltransferase
Nicotinamide-nucleotide amidohydrolase PncC
50S ribosomal protein L28
50S ribosomal protein L33 2
hypothetical protein
Signal recognition particle receptor FtsY
Ribosomal protein S12 methylthiotransferase RimO
DNA-binding protein HU
Integration host factor subunit alpha
ATPase RavA
hypothetical protein
hypothetical protein
von Willebrand factor type A domain protein
von Willebrand factor type A domain protein
photosystem I assembly protein Ycf3
hypothetical protein
Tetratricopeptide repeat protein
hypothetical protein
universal stress protein UspE
Tetratricopeptide repeat protein
DNA gyrase subunit A
ATP-dependent Clp protease ATP-binding subunit ClpC
Chaperone protein HtpG
NTE family protein RssA

Bxylanisolvens_NLAE-zl-P727
Bxylanisolvens_NLAE-zl-P736
Bxylanisolvens_NLAE-zl-P393
Bxylanisolvens_NLAE-zl-P352
Bxylanisolvens_NLAE-zl-P732
Bacteroides_CL03T12C04
Bacteroides_1_1_30
Bacteroides_xylanisolvens_XB1A
Bacteroides_2_1_22
Bacteroides_D1
Bacteroides_xylanisolvens_SD_CC_1b
Bacteroides_D22
Bacteroides_ovatus_SD_CC_2a
Bxylanisolvens_NLAE-zl-C29
Bxylanisolvens_NLAE-zl-H194
Bxylanisolvens_NLAE-zl-G310
Bxylanisolvens_NLAE-zl-G421
Bxylanisolvens_NLAE-zl-C182
Bxylanisolvens_NLAE-zl-C339
Bacteroides_2_2_4
Bacteroides_D2
Bovatus_NLAE-zl-C11
Bovatus_NLAE-zl-C34
Bacteroides_ovatus_CL03T12C18
Bovatus_NLAE-zl-H59
Bxylanisolvens_NLAE-zl-H73
Bacteroides_ovatus_ATCC_8483
Bacteroides_ovatus_3_8_47FAA
Bacteroides_ovatus_3_1_23
Bacteroides_ovatus_SD_CMC_3f
Bovatus_NLAE-zl-H304
Bovatus_NLAE-zl-H361

## B. xylanisolvens non-PUL LGT Event 2



species

site-specific tyrosine recombinase XerD
hypothetical protein
hypothetical protein
Replicative DNA helicase
integration host factor subunit alpha
hypothetical protein
N-acetylmuramoyl-L-alanine amidase
hypothetical protein
hypothetical protein
hypothetical protein
hypothetical protein
hypothetical protein
SusD family protein
TonB-dependent Receptor Plug Domain protein
NADH-dependent butanol dehydrogenase A
hypothetical protein
hypothetical protein
PAP2 superfamily protein
CDP-diacylglycerol--inositol 3-phosphatidyltransferase
GtrA-like protein
Phosphatidylglycerophosphatase A
Myo-inositol-1-phosphate synthase
hypothetical protein
Sensor protein ZraS
Transcriptional regulatory protein ZraR
Outer membrane efflux protein
Autoinducer 2 sensor kinase/phosphatase LuxQ
hypothetical protein
putative ABC transporter permease YknZ
Macrolide export ATP-binding/permease protein MacB
Macrolide export ATP-binding/permease protein MacB
hypothetical protein
ABC transporter ATP-binding protein YtrE
macrolide transporter subunit MacA
NADH dehydrogenase-like protein
Exo-glucosaminidase LytG precursor
Cytidine deaminase
Bifunctional transcriptional activator/DNA repair enzyme AdaA
Inner membrane protein YkgB
Mercuric reductase
putative ABC transporter ATP-binding protein
NADH pyrophosphatase
hypothetical protein
Phosphoglucomutase
Dipeptidase
Glycyl-glycine endopeptidase ALE-1 precursor
hypothetical protein
TonB-dependent Receptor Plug Domain protein
SusD family protein
hypothetical protein
KWG Leptospira
Alanine dehydrogenase
Nitronate monooxygenase
hypothetical protein
hypothetical protein
hypothetical protein
ECF RNA polymerase sigma factor SigW
putative nicotinate-nucleotide pyrophosphorylase [carboxylating]
hypothetical protein
Ribosomal RNA large subunit methyltransferase H

**species**

Bx
Bo

Bxylanisolvens_NLAE-zl-P727
Bxylanisolvens_NLAE-zl-P736
Bxylanisolvens_NLAE-zl-P393
Bxylanisolvens_NLAE-zl-P352
Bxylanisolvens_NLAE-zl-P732
Bacteroides_CL03T12C04
Bacteroides_1_1_30
Bacteroides_xylanisolvens_XB1A
Bacteroides_2_1_22
Bacteroides_D1
Bacteroides_xylanisolvens_SD_CC_1b
Bacteroides_ovatus_SD_CC_2a
Bacteroides_D22
Bxylanisolvens_NLAE-zl-C29
Bxylanisolvens_NLAE-zl-H194
Bxylanisolvens_NLAE-zl-G310
Bxylanisolvens_NLAE-zl-G421
Bxylanisolvens_NLAE-zl-C182
Bxylanisolvens_NLAE-zl-C339
Bacteroides_2_2_4
Bacteroides_D2
Bovatus_NLAE-zl-C11
Bovatus_NLAE-zl-C34
Bacteroides_ovatus_CL03T12C18
Bovatus_NLAE-zl-H59
Bovatus_NLAE-zl-H73
Bacteroides_ovatus_ATCC_8483
Bacteroides_CL02T12C04
Bacteroides_ovatus_3_8_47FAA
Bacteroides_3_1_23
Bacteroides_ovatus_SD_CMC_3f
Bovatus_NLAE-zl-H304
Bovatus_NLAE-zl-H361

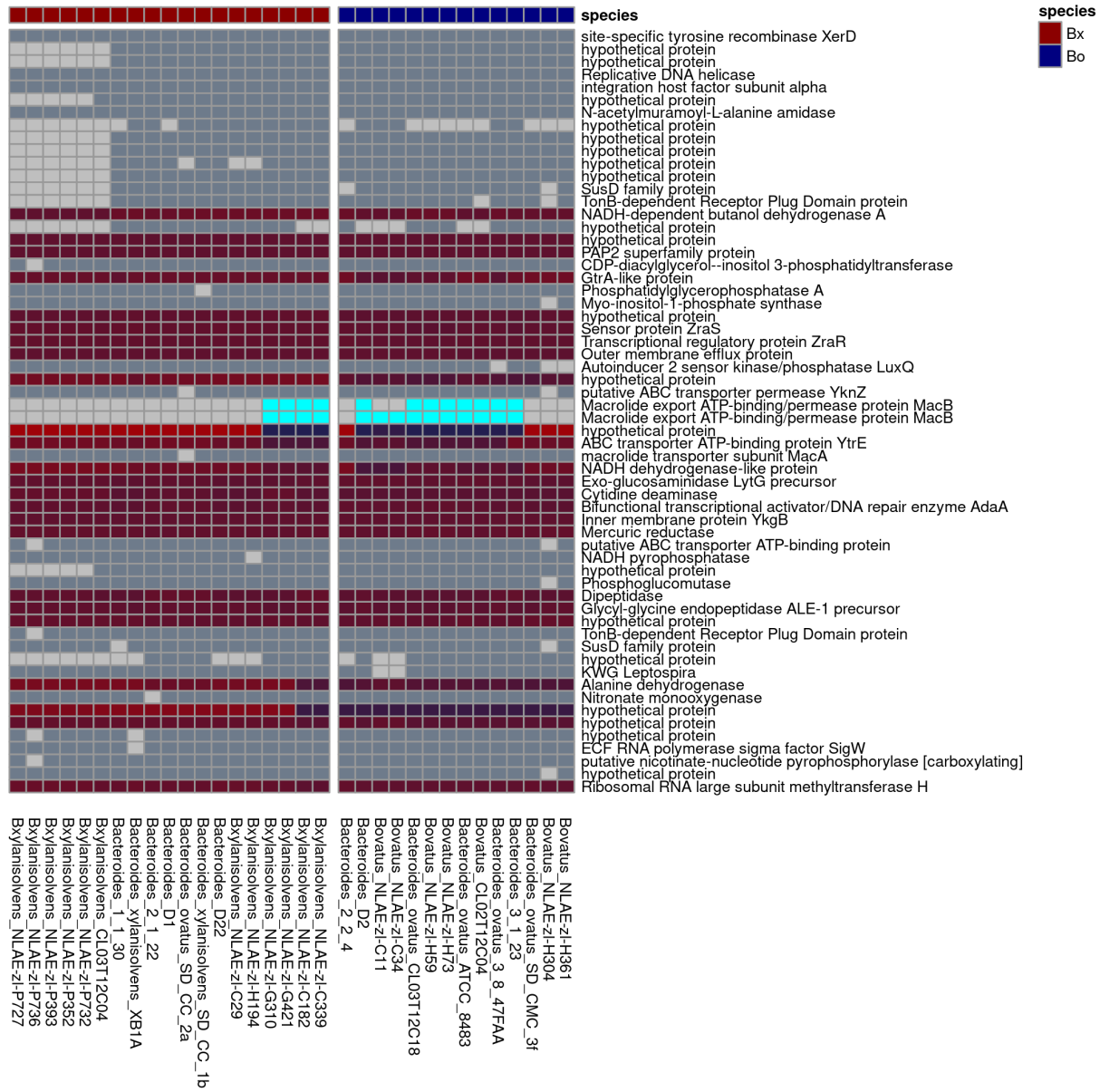**B. xylanisolvens non-PUL LGT Event 3**

**B. xylanisolvens non-PUL LGT Event 4**



species
- Bx
- Bo

Rows (top to bottom):
- hypothetical protein
- Ribosomal RNA small subunit methyltransferase G
- putative metallo-hydrolase
- Glycine dehydrogenase (decarboxylating)
- Putative NAD(P)H nitroreductase YdjA
- HpaII restriction endonuclease
- hypothetical protein
- IMPACT family member YigZ
- hypothetical protein
- D-3-phosphoglycerate dehydrogenase
- Phosphoserine aminotransferase
- ATP-dependent RNA helicase DbpA
- hypothetical protein
- Na(+)-translocating NADH-quinone reductase subunit F
- Na(+)-translocating NADH-quinone reductase subunit E
- Na(+)-translocating NADH-quinone reductase subunit D
- Na(+)-translocating NADH-quinone reductase subunit C
- Na(+)-translocating NADH-quinone reductase subunit B
- Na(+)-translocating NADH-quinone reductase subunit A
- Aminopeptidase E
- Acyltransferase family protein
- Surface antigen
- Outer membrane efflux protein
- Chain length determinant protein
- O-Antigen ligase
- Poly-beta-1%2C6-N-acetyl-D-glucosamine synthase
- N-glycosyltransferase
- Mannosylfructose-phosphate synthase
- N-acetylmuramoyl-L-alanine amidase
- hypothetical protein
- hypothetical protein
- putative AAA-ATPase
- hypothetical protein
- DNA primase
- hypothetical protein
- hypothetical protein
- hypothetical protein
- Teichuronic acid biosynthesis protein TuaB
- hypothetical protein
- hypothetical protein
- putative poly(glycerol-phosphate) alpha-glucosyltransferase
- N-glycosyltransferase
- hypothetical protein
- Poly-beta-1%2C6-N-acetyl-D-glucosamine synthase
- Autoinducer 2 sensor kinase/phosphatase LuxQ
- Endonuclease/Exonuclease/phosphatase family protein

Column labels (left to right):
- Bxylanisolvens_NLAE-zl-P727
- Bxylanisolvens_NLAE-zl-P736
- Bxylanisolvens_NLAE-zl-P393
- Bxylanisolvens_NLAE-zl-P352
- Bxylanisolvens_NLAE-zl-P732
- Bxylanisolvens_CL03T12C04
- Bacteroides_1_1_30
- Bacteroides_xylanisolvens_XB1A
- Bacteroides_2_1_22
- Bacteroides_D1
- Bacteroides_ovatus_SD_CC_2a
- Bacteroides_xylanisolvens_SD_CC_1b
- Bacteroides_D22
- Bxylanisolvens_NLAE-zl-C29
- Bxylanisolvens_NLAE-zl-H194
- Bxylanisolvens_NLAE-zl-G310
- Bxylanisolvens_NLAE-zl-G421
- Bxylanisolvens_NLAE-zl-C182
- Bxylanisolvens_NLAE-zl-C339
- Bacteroides_2_2_4
- Bacteroides_D2
- Bovatus_NLAE-zl-C11
- Bovatus_NLAE-zl-C34
- Bacteroides_ovatus_CL03T12C18
- Bovatus_NLAE-zl-H59
- Bovatus_NLAE-zl-H73
- Bacteroides_ovatus_ATCC_8483
- Bacteroides_CL02T12C04
- Bacteroides_ovatus_3_8_47FAA
- Bacteroides_3_1_23
- Bacteroides_ovatus_SD_CMC_3f
- Bovatus_NLAE-zl-H304
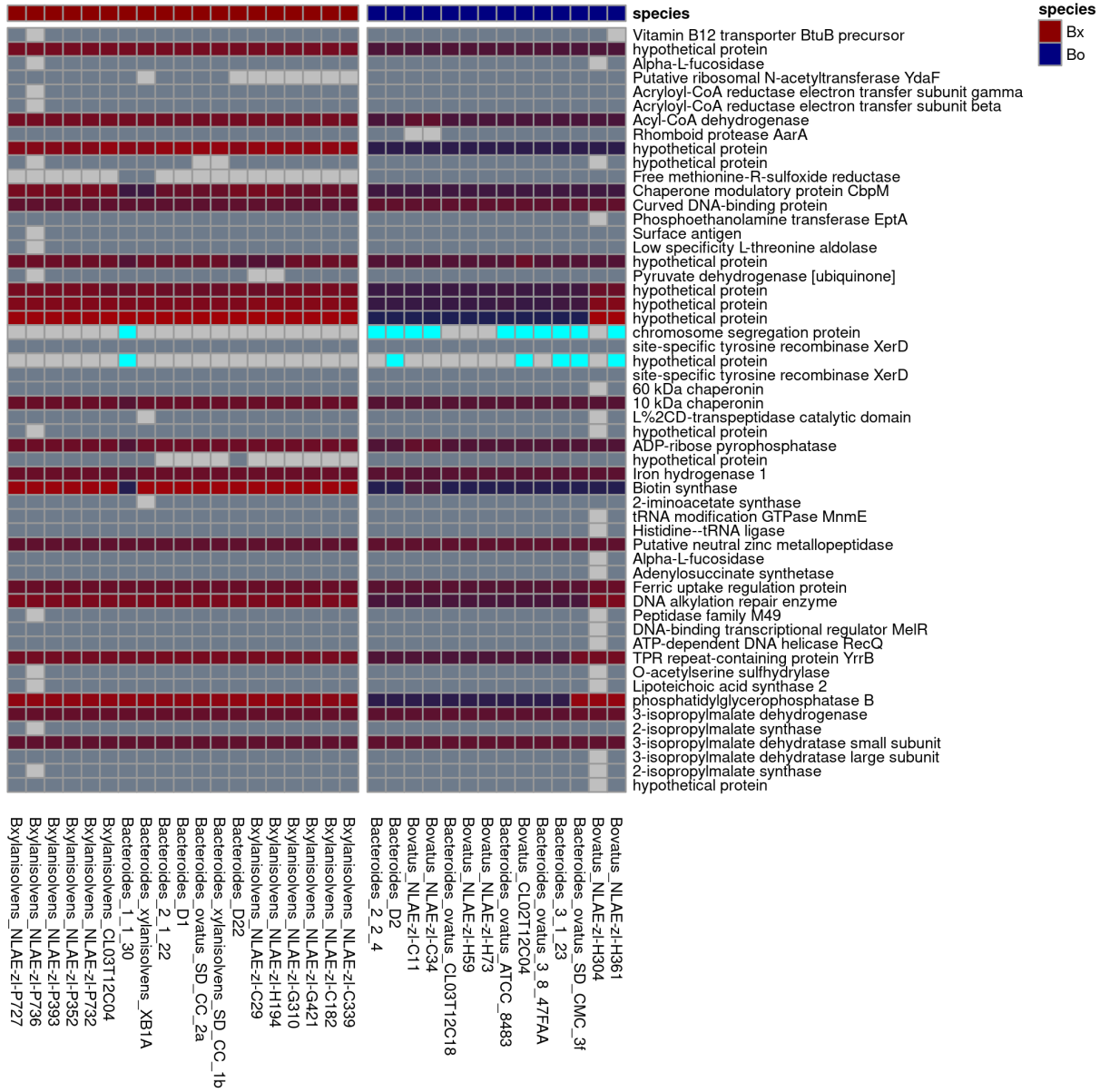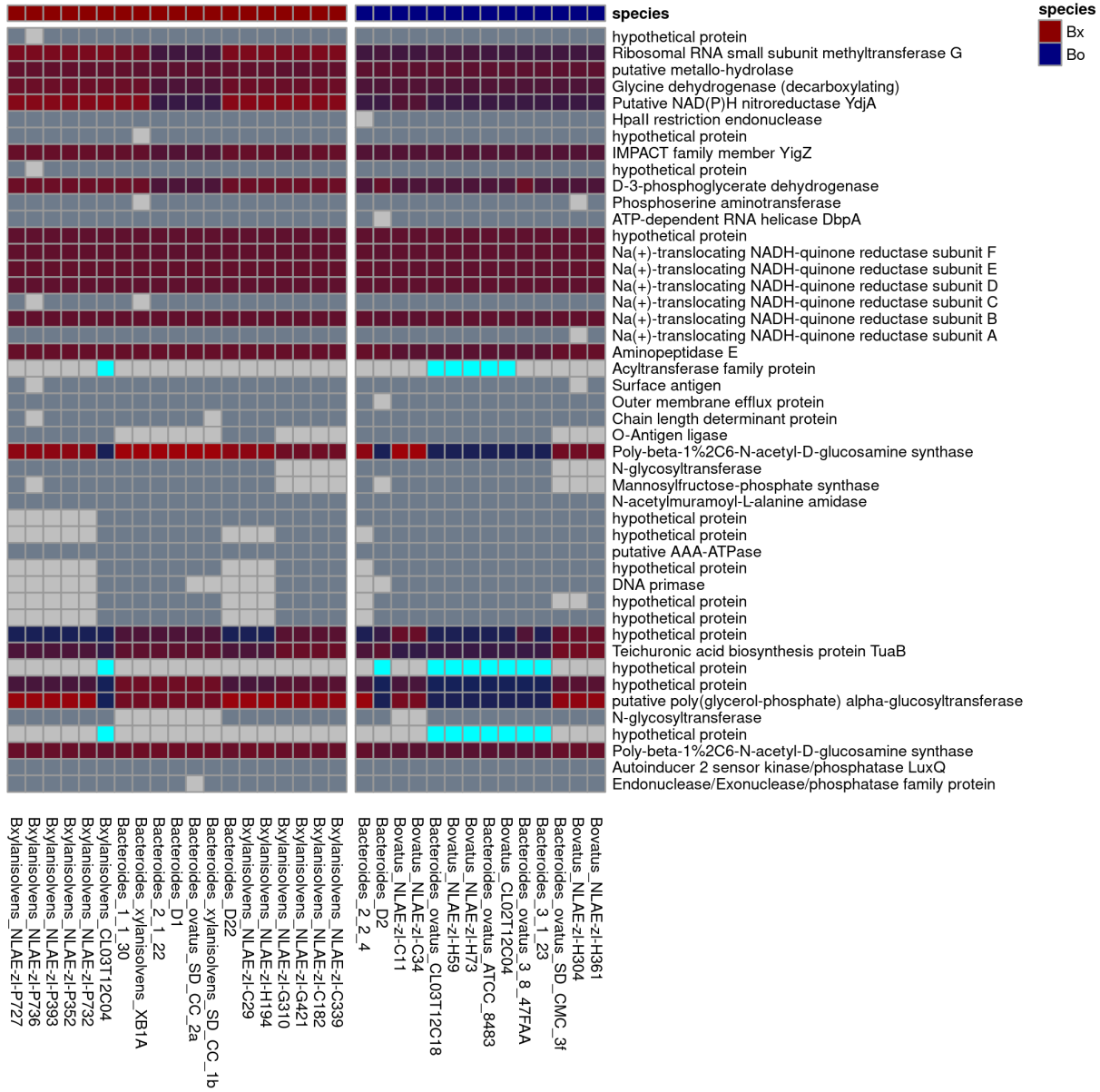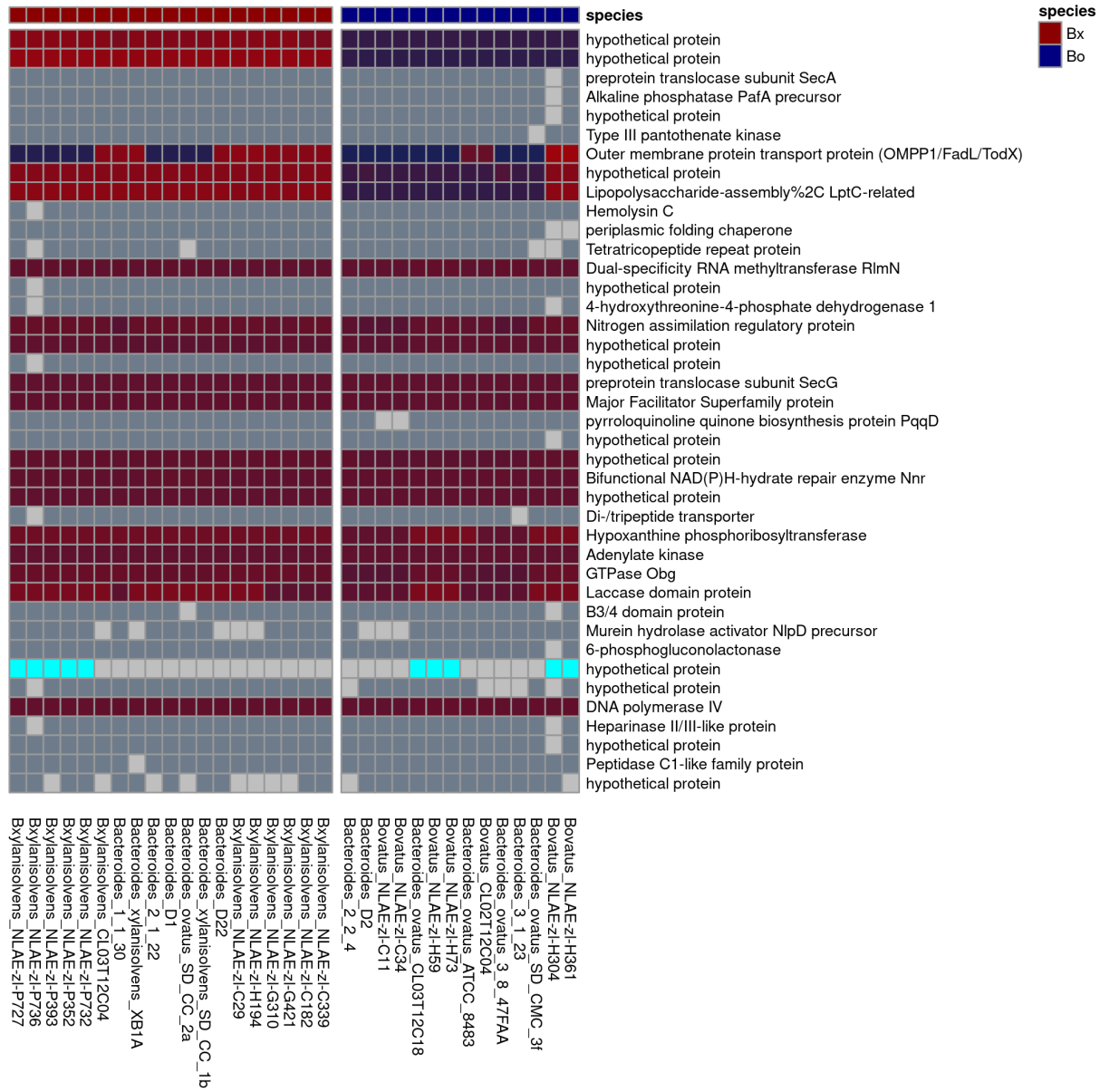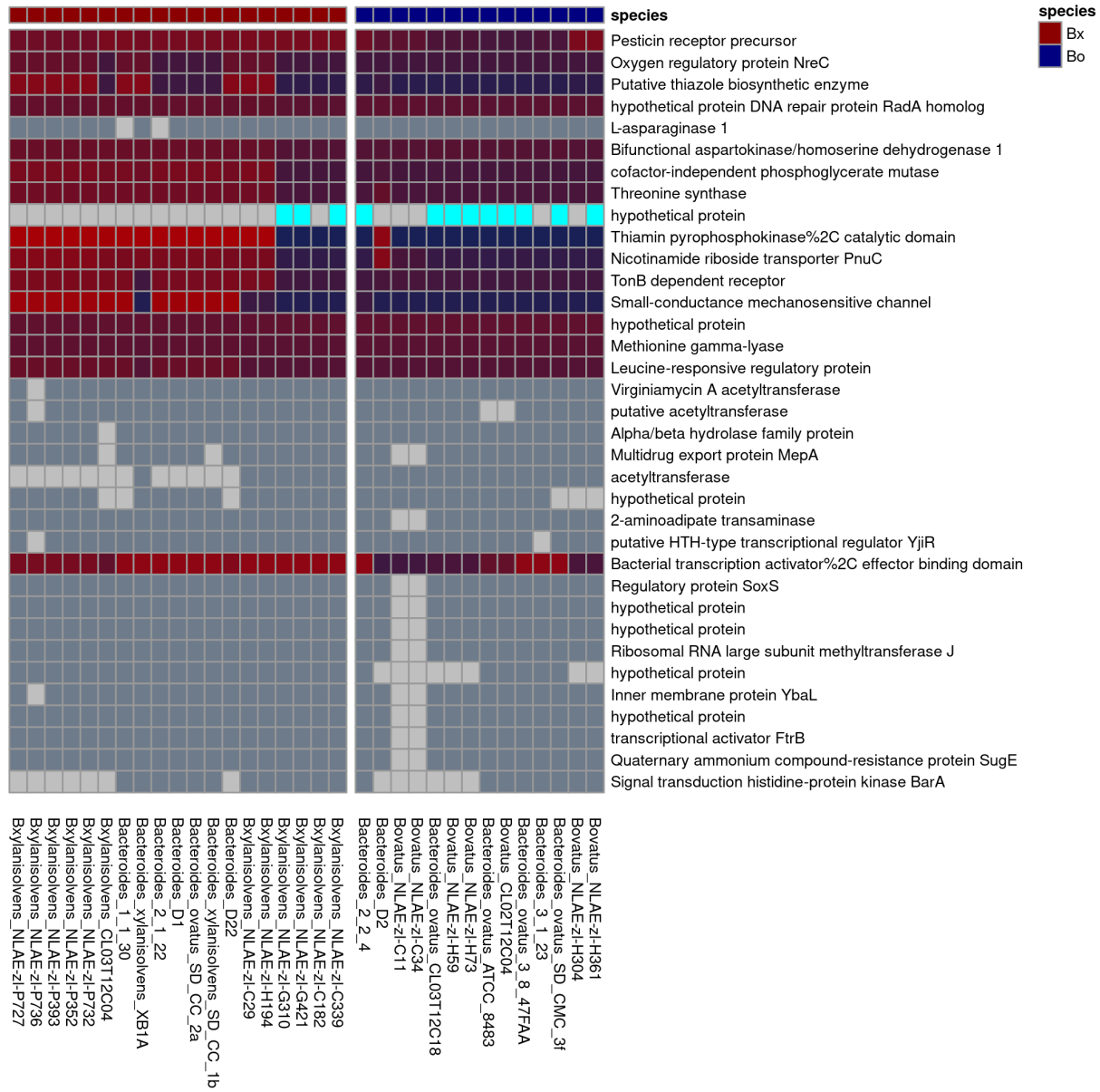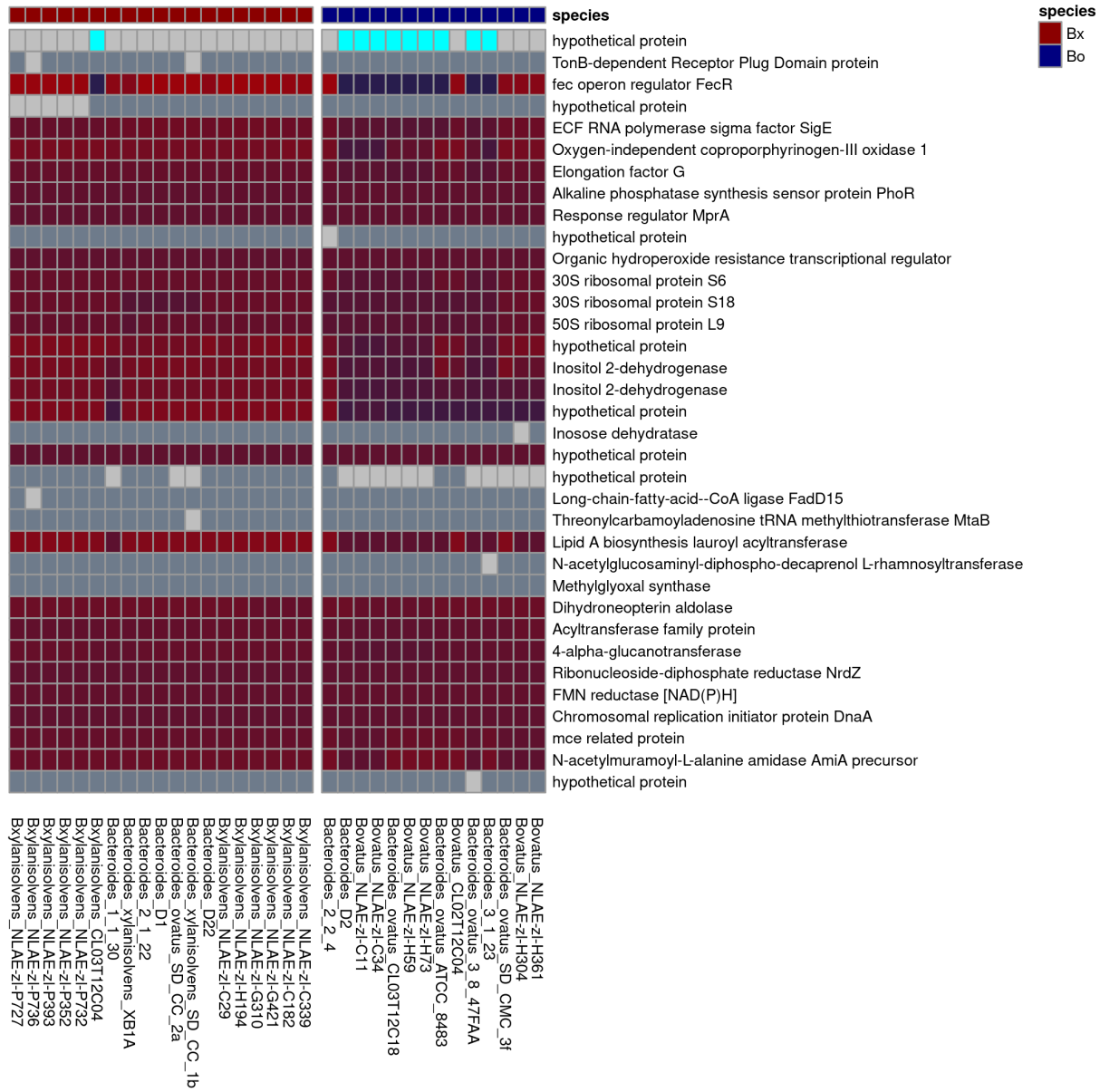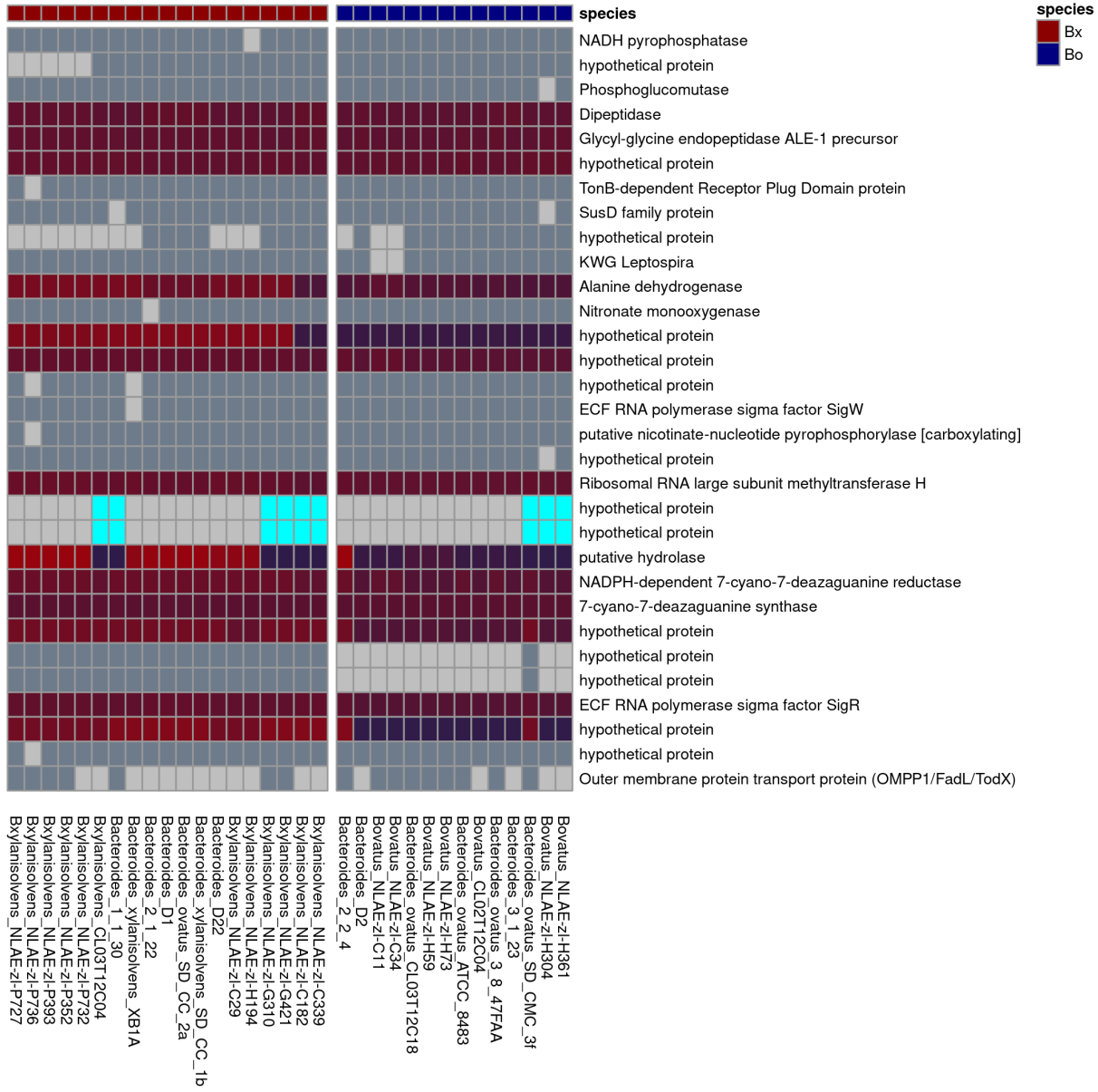- Bovatus_NLAE-zl-H361

## B. xylanisolvens non-PUL LGT Event 5
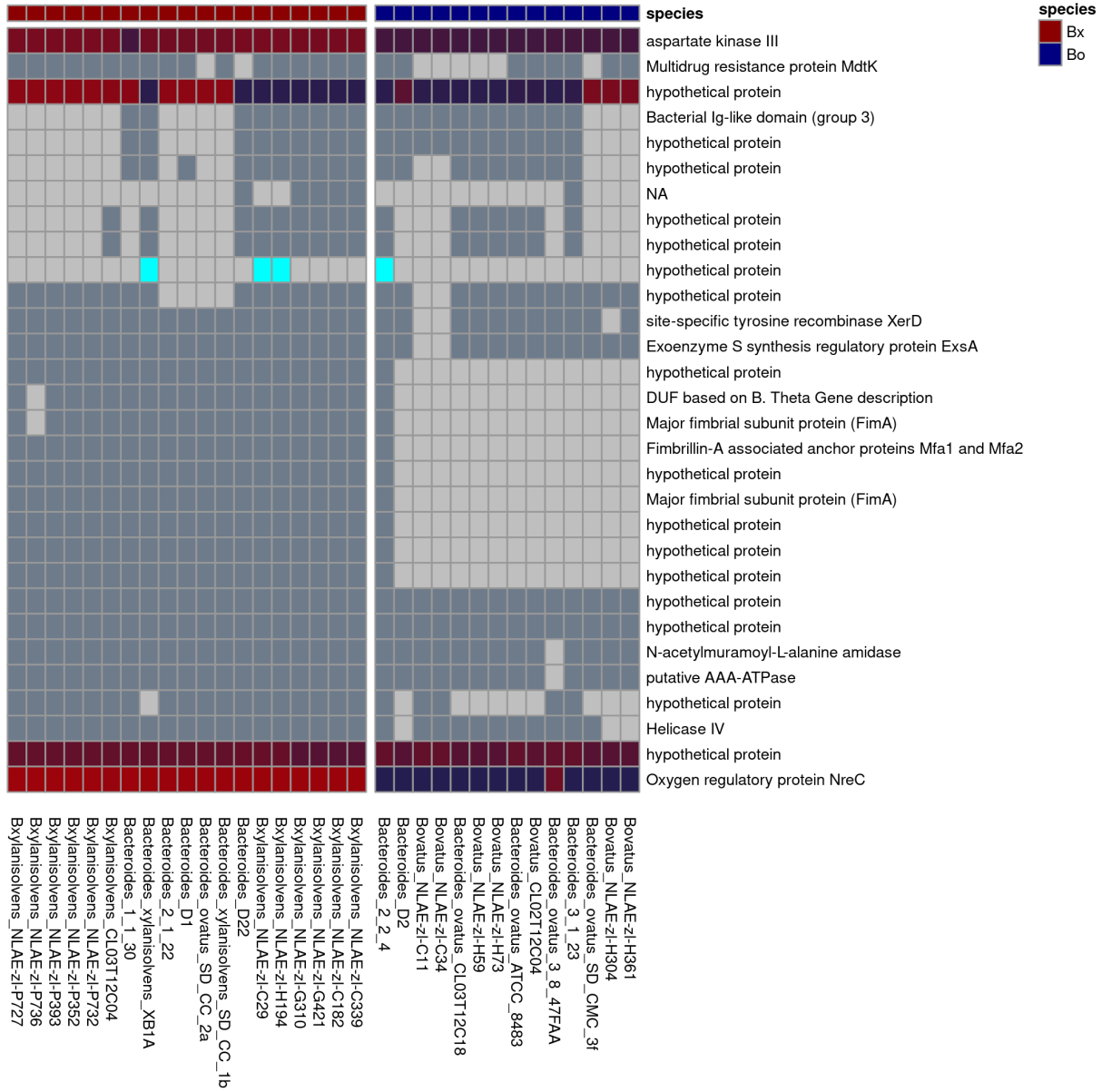
**B. xylanisolvens non-PUL LGT Event 6**

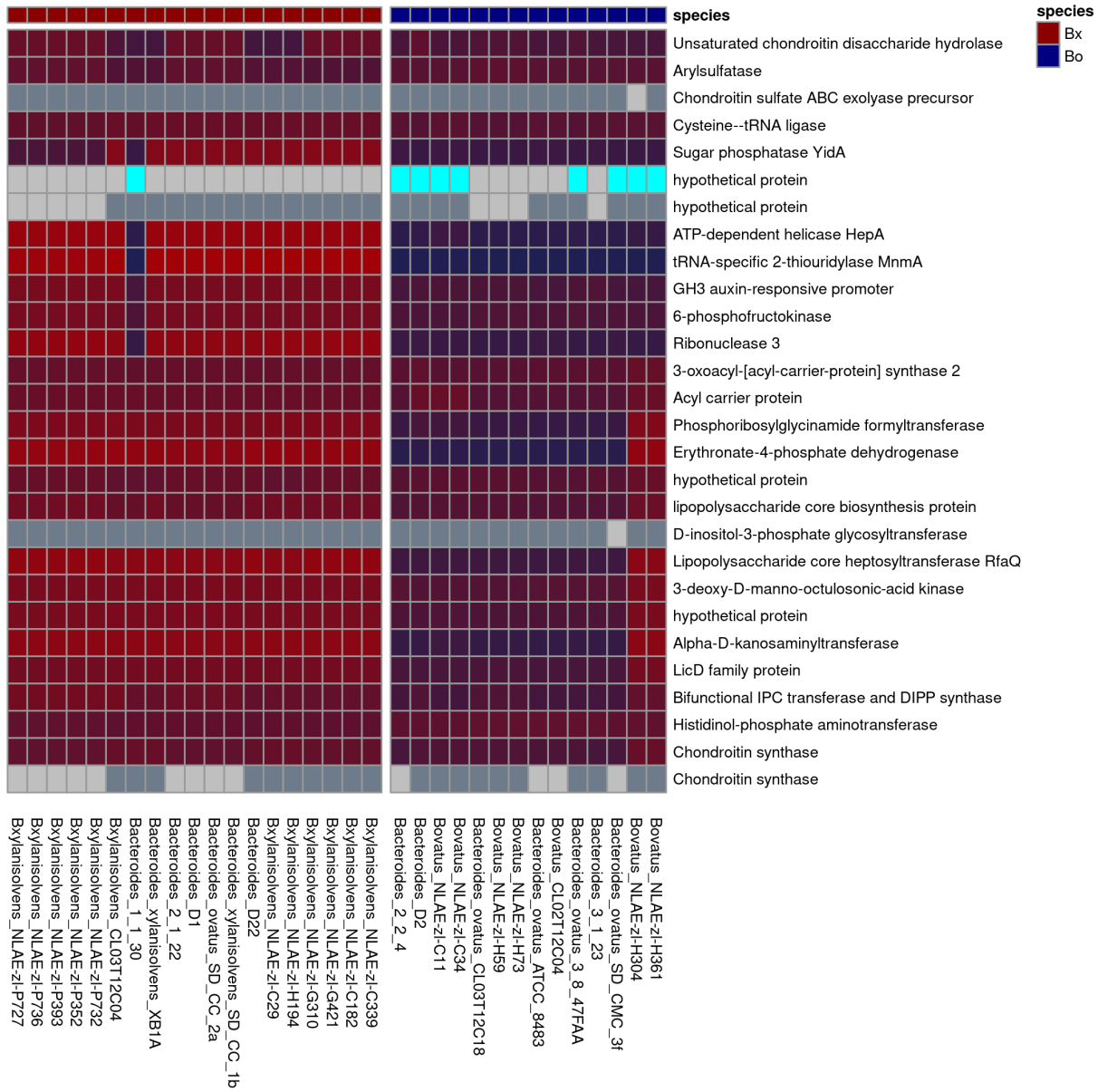**B. xylanisolvens non-PUL LGT Event 7**
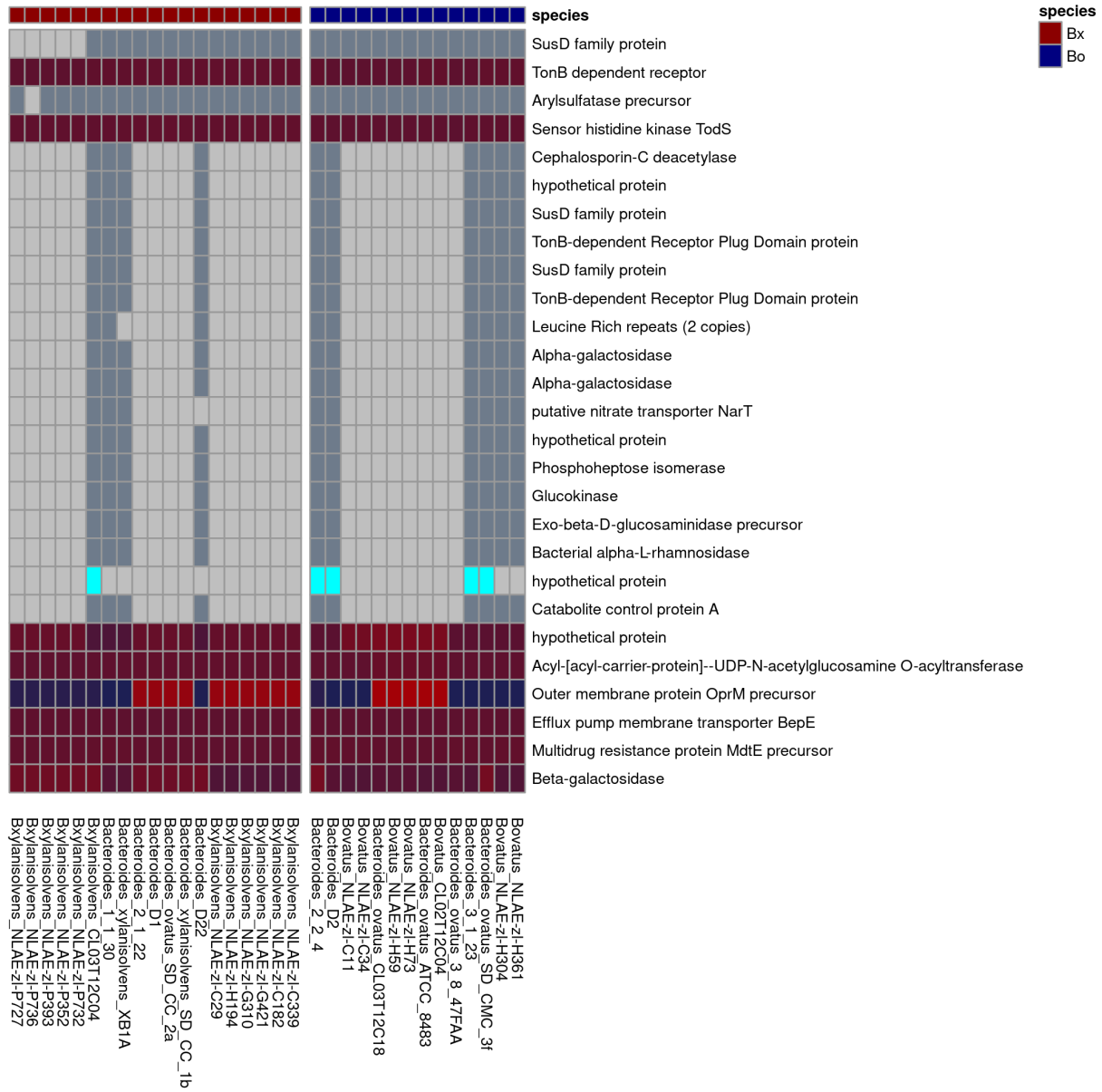
## B. xylanisolvens non-PUL LGT Event 8

## B. xylanisolvens non-PUL LGT Event 9

## B. xylanisolvens non-PUL LGT Event 10

## B. xylanisolvens non-PUL LGT Event 11
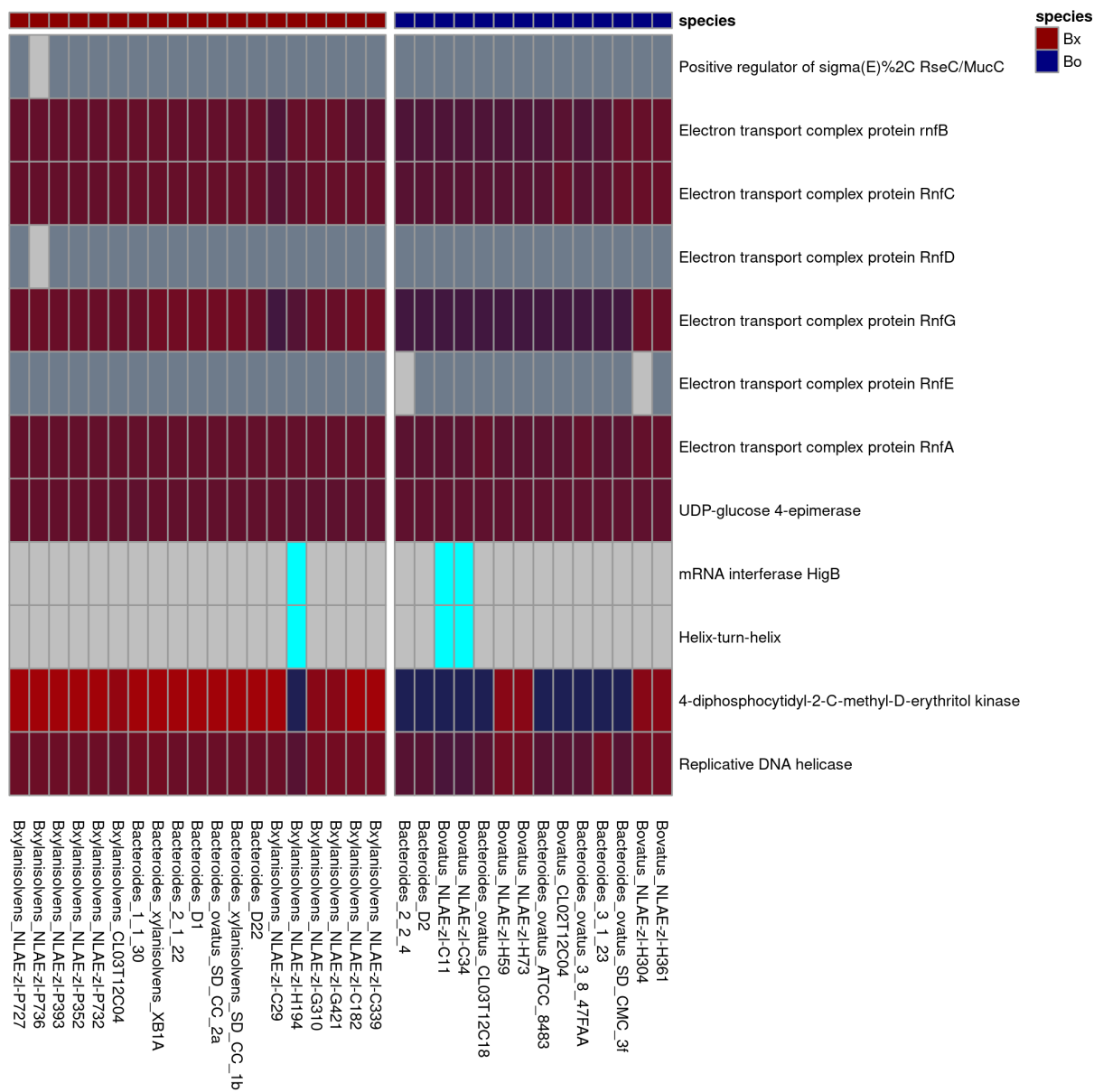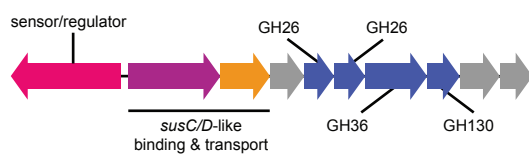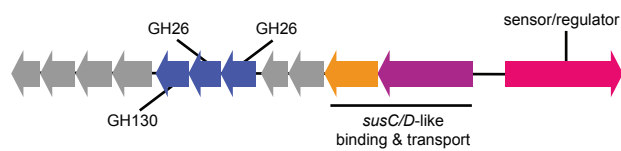


28

**xylanisolvens non-PUL LGT Event 12**

**B. xylanisolvens non-PUL LGT Event 13**

Figure S9



**A.** GalM, GluM PUL-A

**B.** PUL-B expression during GalM growth