1    **Unravelling the specificity and mechanism of sialic acid recognition by the gut**
2    **symbiont *Ruminococcus gnavus***
3

4    C David Owen[1, #,†], Louise E Tailford[2,#], Serena Monaco[3], Tanja Šuligoj[2], Laura Vaux[2],

5    Romane Lallement[2], Zahra Khedri[4], Hai Yu[5], Karine Lecointe[2], John Walshaw[2,6], Sandra

6    Tribolo[2],  Marc Horrex[2], Andrew Bell[2], Xi Chen[5], Gary L Taylor[1],  Ajit Varki[4], Jesus Angulo[3]

7    and Nathalie Juge[2,*].

8    [1]Biomolecular Sciences Building, University of St Andrews, KY16 9ST, UK

9    [2]The Gut Health and Food Safety Programme, Quadram Institute Bioscience, Norwich

10   Research Park, Norwich, NR4 7UA, UK

11   [3]School of Pharmacy, University of East Anglia, Norwich NR4 7TJ, UK

12   [4]Glycobiology Research and Training Center (GRTC), Departments of Medicine and Cellular

13   & Molecular Medicine, UC San Diego, La Jolla, CA, 92093-0687, USA

14   [5]Department of Chemistry, University of California-Davis, Davis, CA, 95616, USA

15   [6]School of Computing Sciences, University of East Anglia, Norwich NR4 7TJ, UK

16   [†]Present address: Diamond Light Source Ltd, Rutherford Appleton Laboratory, Didcot, OX11

17   0FA, UK

18

19   [#]contributed equally to the work

20

21   *Corresponding author. Mailing address: Quadram Institute Bioscience, Norwich Research

22   Park, NR4 7UA Norwich, UK. Phone: +44 (0)1603255068. Fax: +44 (0)1603507723. E-mail:

23   nathalie.juge@quadram.ac.uk

24

25

26

27   **Abstract**

28

29   *Ruminococcus gnavus* is a human gut symbiont which ability to degrade mucins is mediated

30   by an intramolecular *trans*-sialidase (*Rg*NanH). *Rg*NanH comprises a GH33 catalytic domain

31   and a sialic acid binding carbohydrate binding module (CBM40). Here we used glycan

32   arrays, STD NMR, X-ray crystallography, mutagenesis, and binding assays to determine the

33   structure and function of *Rg*NanH_CBM40 (*Rg*CBM40). *Rg*CBM40 displays the canonical

34   CBM40 β-sandwich fold and broad specificity towards sialoglycans with millimolar binding

35   affinity towards α2,3- or α2,6-sialyllactose. *Rg*CBM40 binds to mucus produced by goblet

36   cells and to purified mucins, providing direct evidence for a CBM40 as a novel bacterial

37   mucus adhesin. Bioinformatics data show that *Rg*CBM40 canonical type domains are

38   widespread among Firmicutes. Furthermore, binding of *R. gnavus* ATCC 29149 to intestinal

39   mucus is sialic acid mediated. Together, this study reveals novel features of CBMs which

40   may contribute to the biogeography of symbiotic bacteria in the gut.

41

42

43

**Introduction**

The human gut microbiota encompasses a complex community of bacterial species which play a critical role in human health, through their contribution to e.g. polysaccharide digestion, immune system development, pathogen defence[1]. Microbiota composition varies longitudinally along the gastrointestinal (GI) tract but also transversally from the lumen to the mucosa[1,2]. Most gut bacteria reside in the colon, reaching $10^{11}$ to $10^{12}$ cells per gram, where they compete for dietary and host glycans[3,4]. A dysbiosis of the gut microbiota is associated with intestinal diseases, including cancers, infections, and inflammatory bowel diseases[5-8], underscoring the importance of understanding these host-microbe interactions in order to devise novel treatment strategies.

Several factors influence the biogeography of symbiotic bacteria within the gut, including the gradient and availability of glycans within discrete physical niches[2,3]. The mucus layer covering the GI tract is at the interface between the gut microbiota and the host[5]. In the colon, the mucus layer is divided into a loose outer layer providing a habitat to commensal bacteria and an inner layer adhering to the epithelium and providing protection from bacterial invasion[5]. The outer mucus layer hosts a distinct intestinal microbial niche[9]. The intestinal mucus layers are built around large highly glycosylated gel-forming mucin MUC2 (Muc2 in mouse) secreted by goblet cells[10]. The glycan structures present in mucins are diverse and complex and consist of four core mucin-type O-glycans containing *N*-acetylgalactosamine (GalNAc), galactose (Gal) and *N*-acetylglucosamine (GlcNAc). Mucin O-glycosylation starts with the attachment of GalNAc residues to the hydroxyl group of Ser and Thr of the protein backbone to form the Tn antigen (GalNAcα1-Ser/Thr). This glycan is then elongated into core 1 (Galβ1-3GalNAcα1-Ser/Thr, also known as Thomsen Friedenreich-TF- or T-antigen), core 2 (Galβ1-3(GlcNAcβ1-6)GalNAcα1-Ser/Thr), core 3 (GlcNAcβ1-3GalNAcα1-Ser/Thr) or core 4 (GlcNAcβ1-3(GlcNAcβ1-6)GalNAcα1-Ser/Thr)[11]. Core 3-derived *O*-glycans are important components of human colonic mucin-type *O*-glycans[12]. These core structures are further elongated by the addition of other carbohydrates (e.g. N-acetyllactosamine, LacNAc) and are most commonly terminated by fucose and sialic acid sugar residues *via* α1–2/3/4 and α2–3/6 linkages, respectively. These oligosaccharide chains provide binding sites and nutrients to the bacteria which have adapted to the mucosal environment[13,14]. Reflecting the structural diversity of mucin glycans and their prime location, commensal and pathogenic microbes have evolved a range of adhesins allowing their interaction with mucus[13,15]. Variation in mucosal carbohydrate availability leads to variations in the composition of the resident microbiota[3,16,17] and may also impact on bacterial tropism along and across the GI tract[18].

79    Sialic acids such as *N*-acetylneuraminic acid (Neu5Ac) and fucose residues in terminating

80    positions on mucin glycan chains are prominent targets for commensal and pathogenic

81    bacteria[19,20]. The ratio of sialic acid to fucose increases along the GI tract, from the ileum to

82    the rectum in humans[21] and an inverse gradient occurs in mice[22]. Furthermore blood group

83    Sd[a]/Cad related epitopes, GalNAcβ1-4(NeuAcα2-3)Gal, increase along the length of the

84    human colon[12]. Over 100 complex oligosaccharides can be identified in mucins from human

85    colonic biopsies, with most mono-, di- or trisialylated[23]. Release of sialic acid by microbial

86    sialidases allows bacteria to access free sialic acid for catabolism, decrypt host ligands for

87    adherence, participate in biofilm formation, modulate immune function by metabolic

88    incorporation, and expose the underlying glycans for further degradation[10,14,19,20]. Sialidases

89    are often associated with additional domains including carbohydrate binding modules

90    (CBMs) such as sialic acid specific CBM40[14,24] and broadly specific CBM32[25]. CBMs can

91    enhance catalytic activity by concentrating the enzymes onto carbohydrate substrates[26] or

92    mediate adherence to host cells[27].

93    *Ruminococcus gnavus* is a prominent member of the gut microbiota of the healthy human

94    gut[28]. *R. gnavus* utilisation of mucin is associated with the expression of an intramolecular

95    *trans*-sialidase (IT-sialidase)[29,30], which is proposed to play a key role in the adaptation of gut

96    bacteria to the mucosal environment by providing 2,7-anhydro-sialic acid as a preferential

97    source of nutrients[31]. The IT-sialidase from *R. gnavus* ATCC 29149 (*Rg*NanH) comprises a

98    catalytic glycoside hydrolase domain, *Rg*GH33 and a carbohydrate binding module,

99    *Rg*CBM40.

100    Here, to gain insights into the role and specificity of sialic acid recognition by *R. gnavus*, we

101    employed glycan microarray, X-ray crystallography, saturation transfer difference nuclear

102    magnetic resonance spectroscopy (STD NMR), isothermal titration calorimetry (ITC),

103    mutational analyses, and cell/tissue binding assays to identify *Rg*CBM40 oligosaccharide

104    binding partners. Prominent ligands were oligosaccharides with terminal sialic acid, including

105    those which are not substrates for *Rg*NanH activity. We propose a novel role for CBM40 in

106    targeting gut bacteria towards sialic acid-rich regions of the GI tract.

107

108

109

110

111

112

**Results**

**RgCBM40 belongs to the CBM40 subfamily**

RgCBM40 crystallised as a dimer, adopting the canonical CBM40 β-sandwich fold with six antiparallel strands on the convex face and five on the concave face (**Fig. 1a**, for data collection and refinement statistics, see **Table 1**). Electron density was observed for all RgCBM40 residues present in the construct (50–237). The sialic acid binding site is on the concave face at the dimer interface (**Fig. 1b**), however size exclusion chromatography with multi angle light scattering (SEC-MALS) indicated that the full-length protein, RgNanH, is monomeric in solution (**Fig. 1c**). The macromolecular architecture of RgCBM40 is conserved among members of the CBM40 family (**Supplementary Fig. 1**), with the exception of Vibrio cholerae CBM40_NanH (VcCBM40_NanH) which is proposed to be part of a separate CBM40 subfamily (**Supplementary Fig. 1h**)[25,32]. Greatest structural homology was observed to MdCBM40 NanL (RMSD: 0.3 Å) from the Macrobdella decora IT-sialidase (**Supplementary Fig. 1e**)[33].

Protein ligand complexes were achieved for both 3'SL and 6'SL (**Fig. 1d and e).** No significant conformational changes were observed in the binding site upon ligand binding. Definitive electron density for the Neu5Ac and galactose residues was observed in the 3'SL and 6'SL complexes. In the 6'SL complex, electron density was also observed for the glucose residues (**Fig. 1e**), with the lactose positioned almost perpendicular to the sialic acid (**Fig. 1e**). Contrastingly, for the 3'SL complex, only partial electron density was observed for the glucose residue in a single monomer (**Fig. 1d**), and the glucose positioning indicates that the lactose points up and away from the binding site, without further interactions with the protein. In the 3'SL complex, the lactose positioning would permit further extensions to the carbohydrate chain as would be present in more complex or anchored glycans, whereas these may be blocked in the 6'SL complex. This would provide a degree of specificity towards sialic acid linkage.

Neu5Ac binds in a chair conformation (**Fig. 1f and g**), mimicking the solution conformation and minimizing the energetic penalty paid upon binding[26]. Notably, the carboxylic acid group of Neu5Ac forms electrostatic interactions with an arginine dyad, Arg204 and Arg128, mimicking the coordination observed in sialidase active sites. The C4 hydroxyl group hydrogen bonds to Lys135 and Glu126, the N-acetyl group sits in a hydrophobic pocket formed by Tyr116 and Ile95. The N-acetyl group nitrogen interacts with both Glu126 and Tyr210. Glu126, Arg128, and Arg204 make extensive interactions with the bound ligand and are conserved in all structurally characterized CBM40 sialic acid binding sites, discounting VcCBM40_NanH[34] (**Supplementary Fig. 2**). The environment of the glycerol side-chain of

149    sialic acid is generally conserved across the canonical CBM40 subfamily with the rear face

150    (C7-H and C9-H groups) residing on a hydrophobic surface formed by Ile95 and Tyr210 in

151    *Rg*CBM40 (**Supplementary Fig. 3a**). Although *Vc*CBM40_NanH shares the CBM40 β-

152    sandwich fold (**Supplementary Fig. 1**), the location, orientation, and constitution of its sialic

153    acid binding site is not conserved (**Supplementary Fig. 2**).

154

155    **Structure-based sequence alignment**

156    CBM40s associated with sialidases fall into two subfamilies, the canonical subfamily

157    exemplified by *Cp*CBM40_NanJ[25] (which also regroups *Rg*CBM40, *Cp*CBM40_NanI[32],

158    *Sp*CBM40_NanA[35], *Sp*CBM40_NanB[36], *Sp*CBM40_NanC[37] and *Md*CBM40 NanL[33]),  and the

159    *Vibrio* subfamily exemplified by *Vc*CBM40_NanH[34]. Considerable sequence divergence

160    between the *Vibrio* and canonical CBM40 types renders satisfactory alignments difficult to

161    produce with standard tools, as also previously reported[32]. Here, by detailed manual

162    inspection, paying particular attention to the limits of secondary structure elements and

163    intervening loops, we produced an alignment of both types of CBM40 sequences showing

164    well-conserved positions along its length, notwithstanding the *Vibrio* insertion (40 residues)

165    near the N-terminus. The pairwise identities between the canonical representatives range

166    from 21–67%, while the maximum canonical versus *Vibrio* identity is 17%, reflecting that

167    CBM40s fall into two distinct groups. This highlighted conserved residues within the

168    canonical subfamily that may be involved in binding affinity and specificity (**Fig. 2**). These

169    include (*Rg*CBM40 numbering): an arginine dyad (Arg204 and Arg128) that interacts with

170    the sialic acid carboxylic acid group, a glutamic acid (Glu126), which hydrogen bonds to the

171    C4 hydroxyl; and a hydrophobic surface, which accommodates the N-acetyl moiety and the

172    hydrophobic face of the glycerol group. Tyr116, Ile95, Tyr210 contribute to the surface of an

173    aromatic:aliphatic:aromatic twisted platform which presents the glycerol hydroxyl groups to

174    solvent[26].

175    **Bioinformatics analyses**

176    To gain further insights into the phylotypic distribution of the CBM40 domains within bacterial

177    genomes, we performed a database search using pHMMs derived from our alignment as

178    queries (canonical and *Vibrio*-type together, referred to as 'combined'; canonical only; *Vibrio*-

179    type only) as well as Pfam models, "Sialidase(NTD)", "Laminin_G_3", and "Sial-lect-inser"

180    (see **Methods** and **Supplementary Methods**). Our combined model successfully identified

181    99.9% of the CBM40 domains matched by the individual type CBM40 models (over 16,000

182    domain hits in the whole database of around 67,000 genomes). Further analysis of the data

183    (see **Supplementary Methods**) led to the identification of 51 nonredundant sequences

184    (**Supplementary Fig. 4**). Of these, the canonical CBM40 domains occurred in Firmicutes

185   with 40 sequences, representing 18 genera or pseudogenera, divided between classes

186   Bacilli and Clostridia, as well as Erysipelotrichi and an unclassified member of the

187   Firmicutes; and two sequences in Actinobacteria. The Vibrio type occurred only in

188   Gammaproteobacteria, represented by 8 sequences in five genera. The separation between

189   the Vibrio-type sequences and canonical CBM40 sequences across bacterial genomes was

190   also apparent from a tree representation constructed using a simple distance-based model

191   and neighbour-joining (**Fig. 3**). This dichotomy was fully supported by bootstrap analysis of

192   1,000 replicates. There was no evidence for any intermediate or other CBM40 types. Only

193   one sequence from *Actinobacillus muris* containing a canonical CBM40 (confirmed by

194   pHMMs and conserved binding residues) was shown to be part of a Gammaproteobacteria

195   clade (all other members Vibrio type) as supported by 79% of bootstraps. Further studies

196   may indicate whether this domain is the closest to an inferred common ancestor of the

197   canonical and Vibrio CBM40 types. The results for co-incidence of sialidase domains clearly

198   indicated an association with CBM40s in this set of nonredundant sequences: we detected a

199   sialidase domain in 92% of canonical-type CBM40 and in all Vibrio type CBM40

200   representatives.

201

### *Rg*CBM40 preferentially binds α2,3 linked sialosides

203   To further explore *Rg*NanH ligand specificity, *Rg*CBM40 and inactive mutant *Rg*GH33

204   D282A, were tested for binding to various sialoglycans, using a slide microarray[38,39]. This

205   sialoglycan microarray presents over 60 synthetically recreated naturally-occurring

206   oligosaccharide structures with diverse sialic acid forms, glycosidic linkages, and underlying

207   glycans, representing a broad range of such targets[38,39]. Both recombinant proteins

208   exclusively bound to glycans terminated with sialic acids (**Fig. 4**). They also showed distinct

209   specificities. *Rg*CBM40 bound to terminal Neu5Ac, Neu5Gc, Neu5,9Ac$_2$ and 2-keto-3

210   deoxynonulosonic acid (Kdn) attached with α2-3, α2-6 and α2-8 linkages (**Fig. 4**). In

211   contrast, *Rg*GH33 D282A interacted weakly with a narrow spectrum of sialoglycans, mainly

212   α2-3-Neu5Ac-containing glycans, primarily Neu5Acα3LacNAcβ (3'SLN),

213   Neu5Acα3Galβ3GlcNAcβ, Neu5Acα3Galβ3GalNAc (STF), Neu5Acα3Lacβ (GM3), and

214   Neu5Acα3Galβ3GalNAcβ3Lac (**Fig. 4**). Noticeably, *Rg*GH33 D282A recognized some of the

215   α2-3-linked sialoglycans but not any α2-6- or α2-8-linked ones, in line with its substrate

216   specificity[30]. In marked contrast, every α2-3-linked sialyl oligosaccharide present on the

217   array could be bound by *Rg*CBM40. *Rg*CBM40 showed a preference for terminal Neu5Ac

218   over Neu5Gc, and for α2-3>>α2-6>α2-8 linkages. *Rg*CBM40 bound generally more strongly

219   to glycans containing LacNAc and Lac. *Rg*CBM40 could bind Neu5Ac linked Lac with α2-3

220   and α2-6 linkage, albeit to a lesser degree, whereas binding to Neu5Ac linked LacNAc was

221   α2-3-specific. Due to the glycan orientation introduced by the α2-6-sialic acid linkage the

222   6'SL glucose residue is close to the protein surface (**Fig. 1e**). Therefore, α2-6-linked LacNAc

223   *N*-acetyl group may be blocked by protein residues, whereas the α2-3 linked glycan would

224   be more solvent exposed. The highest binding was to Neu5,9Ac$_2$α3GalβR1. Interestingly,

225   *Rg*CBM40 bound to Neu5Gcα3Galβ3GalNAcβR1 (Neu5Gc-TF) and

226   Neu5Gc9Acα3Galβ3GalNAcβR1 (Neu5Gc9Ac-TF) although with 5–10 fold less intensity, but

227   it could not bind to the same ligands with the αR1 linkage. *Rg*CBM40 bound to α2-3-

228   sialylated Lewis X (3'SLX, both Neu5Ac and Neu5Gc forms, although Neu5Ac was

229   preferred). Sulfation of the 6 position of GlcNAc in 3'SLX (both Neu5Ac and Neu5Gc)

230   improved binding of the protein (**Fig. 4**).

231   To validate some of the glycan array data, we used STD NMR spectroscopy[40,41] against a

232   range of sialylated ligands. Since the highest STD intensities correlate with the closest

233   ligand-protein contacts in the bound state[42], STD NMR experiments provide important

234   information on the binding epitope of the complexed ligand[43].

235   Here Neu5Ac, Neu5Gc, 2,7-anhydro-Neu5Ac, 3'SL, 6'SL, Neu5Acα3Gal (3'SGal),

236   Neu5Acα6Gal (6'SGal), 3'SLN, Neu5Acα6LacNAc (6'SLN), Neu5Gcα3Lac (3'SLGc),

237   Neu5Gcα6Lac (6'SLGc), Neu5Acα6GalαOC3H6N3 (Neu5Ac-STn),

238   Neu5Gcα6GalαOC3H6N3 (Neu5Gc-STn), and STFαOC3H6N3 were tested as potential

239   ligands for *Rg*CBM40. With the exception of the three monosaccharides, Neu5Ac, Neu5Gc,

240   and 2,7-anhydro-Neu5Ac, binding to *Rg*CBM40 was detected for all di- and tri-saccharides

241   tested. For the latter, the binding epitope mapping was obtained and analyzed as described

242   under Methods. **Fig. 5a** shows the STD NMR spectra of 3'SL and 6'SL, and **Fig. 5b** their

243   binding epitope mapping. The sialic acid ring was found to be the main recognition element

244   and the binding mode was not affected by the nature of the glycosidic linkage (α2-3 or α2-6)

245   of the sialoglycan (**Supplementary Fig. 5**). The same was true for the other Neu5Ac-ending

246   ligands tested (see binding epitope mapping in **Supplementary Fig. 6**). The overall binding

247   epitopes of 3'SL and 6'SL from the STD NMR in solution state are in good agreement with

248   the crystal structures (**Fig. 5**), where the sialic acid is in close contact to the protein surface

249   while the lactose moiety is solvent exposed as suggested from the very low STD intensities

250   observed for the galactose and glucose protons. Very strong STD intensity is observed at

251   the methyl group (**Fig. 5**). This is in excellent agreement with the N-acetyl group sitting in the

252   hydrophobic pocket facing many protein protons (Hδ and Hγ) from the side chains of Ile95,

253   Tyr116, and Tyr210 (**Fig. 1f and 1g**). High intensity on H7 compared to the much lower one

254   on the adjacent H8 agrees with H7 facing the hydrophobic side chains while H8 (**Fig. 5**), in

255   trans-conformation to it, is pointing towards the solvent. Within experimental error, no stark

256   differences were observed in the orientation of the sialic acid ring in the binding pocket of

8

257    *Rg*CBM40. *Rg*CBM40 also showed binding to Neu5Gc-ending oligosaccharides, albeit with

258    a lower strength. **Fig. 5c** shows the binding epitope of 3'SLGc and 6'SLGc (STD spectra are

259    shown in **Supplementary Fig. 6**). Again, sialic acid was the main recognition element of

260    these sialoglycans, but the binding epitope mapping was slightly different, in comparison to

261    those of 3'SL and 6'SL. For the Neu5Gc-ending ligands, stronger STD intensities on H3s

262    and lower ones on H6 were observed, suggesting a small reorientation of the ring around

263    C3, which would expose C6, in order to fit the bulkier hydroxyl group on the acetamide

264    moiety.

265    The affinity of the interaction between *Rg*CBM40 and sialic acid ligands was further

266    assessed by ITC. Both 3'SL and 6'SL bound with similar low affinities, with dissociation

267    constants of 0.57 mM and 1.70 mM, respectively (**Fig. 6a and b, Supplementary Table 1**).

268    This confirms that *Rg*CBM40 is specific for the terminal residue irrespective of the glycosidic

269    linkage but with a slight preference (~ 3 fold) for the 2-3 linkage. Furthermore, it would

270    suggest that the additional binding interactions observed in the crystal structure of the

271    complex between *Rg*CBM40 and 6'SL do not significantly promote binding, also in

272    agreement with the STD NMR results, showing that sialic acid is the main binding epitope in

273    solution. We confirmed that *Rg*CBM40 binds to Neu5Gc-oligosaccharides, albeit with lower

274    affinity, in accordance with the glycan array and STD NMR results. *Rg*CBM40 has a Kd of

275    ~3 mM and >10 mM towards 3'SLGc and 6'SLGc, respectively (**Fig. 6c, Supplementary**

276    **Table 1**). Very weak (~20 mM) interaction was observed between *Rg*CBM40 and Neu5Ac

277    (**Fig. 6d**) or Neu5Gc monosaccharides (**Supplementary Table 1**). The STD NMR

278    experiments were carried out with 1 mM sugar, well below the Kd, which explains why no

279    interaction was observed using this approach. Thermodynamic analysis showed that the

280    reaction is enthalpy-driven (**Supplementary Table 2**).

281    To further assess the involvement of individual residues we introduced point mutations

282    specifically designed to abrogate CBM binding. Arg128, Arg204, Tyr116, Tyr 210, Glu126

283    and Ile95 were chosen for alanine substitutions. Analysis of the secondary structure by

284    circular dichroism (CD) suggests that the recombinant proteins were correctly folded

285    (**Supplementary Fig. 7**). Binding to Neu5Ac, 3'SL and 6'SL was abolished for the double

286    mutant R128A/R204A as well as all single mutants, with the exception of I95A as shown by

287    ITC (**Supplementary Fig. 8a and b, Supplementary Table 1**). I95A binds 3'SL and 6'SL

288    with a Kd of 1.82 and 1.37 mM, respectively, broadly similar to the binding of the wild type

289    enzyme (**Supplementary Table 1**).This suggests that Ile95 is not an essential component of

290    the hydrophobic pocket or the aromatic:aliphatic:aromatic twisted platform, and that the Tyr

291    residues may compensate for the mutation of Ile95 to Ala. The binding ability of I95A to 3'SL

292    and 6'SL was further confirmed by STD-NMR (**Supplementary Fig. 9**).

293  Taken together, the STD NMR and ITC data confirmed binding of both α2-3 and α2-6 linked
294  sugars and raise questions regarding differences in ligand specificity between the catalytic
295  and carbohydrate binding domains constituting *Rg*NanH. We previously showed that
296  *Rg*NanH is specific for α2-3-linked substrates[30]. To determine the influence of *Rg*CBM40 on
297  the sialidase activity, we compared the enzymatic activity of *Rg*NanH and *Rg*GH33 on a
298  range of sialylated substrates. The reaction was monitored by HPAEC-PAD and showed no
299  difference in catalytic activity on short oligosaccharides 3'SL, 3'SLX (Neu5Ac form) or on
300  large polymeric MUC2 mucins (**Supplementary Fig. 10)**, indicating that, in the conditions
301  tested, *Rg*CBM40 did not potentiate the enzyme activity on these substrates.
302
303  **_Rg_CBM40 is a novel bacterial mucus adhesin**

304  *R. gnavus* ATCC 29149 but not the E1 strain encodes the IT-sialidase required for mucin-
305  degradation[29,30]. Immunogold labeling and western blotting confirmed the presence of
306  *Rg*NanH on *R. gnavus* ATCC 29149 cell-surface but not E1 (**Supplementary Fig. 11a and**
307  **b**). Given the role of *Rg*NanH in *R. gnavus* mucin glycan utilization, the binding of *Rg*CBM40
308  was tested towards a range of mucins with different glycosylation profiles by ELISA. The
309  sialylation level of purified commercial pig gastric mucin (pPGM), mixed and Muc2/MUC2
310  mucins from mice and LS174T human cell line was analyzed by mass spectrometry (MS),
311  revealing that most of the mucins tested contained >8% sialylated structures; pPGM and
312  Muc2 from the colon of wild type C57BL/6 mice contained <2% sialylated structures whereas
313  the level of sialylation of LS174T MUC2 reaches 91% (**Supplementary Table 3**). Highest
314  binding was observed to LS174T MUC2 whereas binding was lowest to pPGM or Muc2 from
315  the colon of wild type mice, which contain low levels of sialylation (**Fig. 7a**). The interaction
316  was dependent on the concentration of *Rg*CBM40 (**Supplementary Fig. 12**). *Rg*CBM40
317  generally bound more strongly to mucins extracted from *C3GnT*[-/-] mice (mutants which lack
318  core 3 β1-3-N-acetylglucosaminyltransferase, C3GnT)[44] than to mucins from wild type mice.
319  Irrespective of the mouse model, the binding of *Rg*CBM40 to Muc2 from the small intestine
320  was higher than from the colon (**Fig. 7a**). The adhesion level correlated well with the level of
321  sialylation between the different mucins tested ($r^2$ = 0.88; **Fig. 7b**). *Rg*CBM40 bound
322  significantly less strongly to MUC2 which has been treated with trifluoroacetic acid (TFA) to
323  remove sialic acid, or with any of the sialidases tested which included the broad-specificity
324  sialidase from *Clostridium perfringens* (*Cp*) and the α2-3-specific sialidases from *Salmonella*
325  *typhimurium* (*St*), *Akkermansia muciniphila* (*Ak*) and *R. gnavus* (*Rg*)**,** confirming the
326  specificity of *Rg*CBM40 for terminal sialic acid (**Fig. 7c**). Consistent with the low affinity of
327  CBM40 for Neu5Ac, this monosaccharide had no effect on adherence of *Rg*CBM40 to mucin
328  (**Fig. 7d**). However, addition of free 3'SL or 6'SL prior to binding significantly decreased

329  adherence of *Rg*CBM40 to MUC2 (**Fig. 7d**). These data indicate that *Rg*CBM40 recognizes

330  sialylated mammalian mucins.

331  Having shown that *Rg*CBM40 can bind to sialylated oligosaccharides and mucins, we tested

332  its ability to bind to mucus from mouse intestinal tissue and human cell lines cells by

333  immunofluorescence (**Fig. 8**). Methacarn fixation allowed preservation of mucus in both

334  tissue sections and cell lines. Strong binding was demonstrated to mucus produced by

335  LS174T which correlated with staining patterns of SNA (a sialic acid specific lectin) and

336  MUC2 (**Fig. 8a**). No staining was observed in negative controls (*Rg*CBM40 free). *Rg*CBM40,

337  Muc2 and lectin staining was also observed in crypts as well as on the epithelial surface of

338  mouse colonic tissue (**Fig. 8b**). In addition, sialidase treatment of mouse colonic sections

339  markedly reduced the binding of *Rg*CBM40 as well as the SNA lectin control (**Fig. 8c**). SNA

340  can outcompete *Rg*CBM40 binding to the mucus layer in mouse colonic tissue sections,

341  further indicating that the binding of *Rg*CBM40 to mucus is sialic acid mediated (**Fig. 8d**).

342  Similar inhibition was observed when using bacterial cells.  *R. gnavus* ATCC 29149 was

343  shown to bind to areas that correlated with mucus staining. This binding was blocked with

344  the addition of SNA (**Fig. 8e**), confirming the importance of sialic acid recognition in *R.*

345  *gnavus* ATCC 29149 binding to mucus.

346

347  **Discussion**

348  Sialic acids are often found capping mammalian glycans and are thus common binding

349  targets of commensal or invading microbes. A wide variety of microorganisms utilize CBM-

350  containing sialidases to process these terminal sialic acid residues. At present CBMs in

351  family 40 are the only known examples to bind sialic acid and are exclusively associated

352  with sialidases (www.cazy.org). The CBM40 from *R. gnavus*, *Rg*CBM40, adopts the

353  characteristic CBM40 β-sandwich fold, previously reported for CBM40s present in *C.*

354  *perfringens*[25, 32], *V. cholerae*[34], *M. decora*[33] as well as *S. pneumoniae*[35,36,37].

355  In their description of *C. perfringens Cp*CBM40_NanJ, Boraston *et al.* pointed out that there

356  appears to be two subfamilies within the CBM40 family, one typified by *Cp*CBM40_NanJ and

357  the other by *V. cholerae Vc*CBM40_NanH[25]. This was further supported by phylogenetic

358  analyses of all CBM40 structurally characterized so far[32]. It is clear that *Vibrio* sp. forms an

359  outlying clade in the family that has very low amino acid sequence identity (<15%) with the

360  main clade[32]. Here, we showed that the separation between the Vibrio-type sequences and

361  canonical CBM40 sequences is also observed across bacterial genomes. Both types adopt

362  a β-sandwich fold, however this is the most common core fold across CBM families[26].

363  *Rg*CBM40 crystal structures, of the canonical type, in complex with sialylated ligands

364    demonstrate shared core binding site residues. In brief, on one side of the sialic acid

365    residue, the carboxylic acid and C4 hydroxyl groups are coordinated by an arginine dyad

366    (Arg128 and Arg204) and a glutamic acid (Glu126) residue, respectively. The importance of

367    the arginine residues was further confirmed by mutational analyses, showing loss of binding

368    of *Rg*CBM40 R204A, *Rg*CBM40 R128A and the double mutant *Rg*CBM40 R128A/R204A to

369    3'SL. The methyl of the N-acetyl moiety and the C-H face of the glycerol moiety reside on a

370    hydrophobic twisted platform surface formed by primarily aromatic residues, of which Tyr116

371    and Tyr210 are essential for binding. Glu126 was also shown to be essential, as predicted

372    given its conservation and interactions with both the *N*-acetyl group N and the C4 hydroxyl of

373    the sialic acid moiety.

374    *Rg*CBM40 showed broad specificity for sialylated oligosaccharides with dissociation

375    constants to 3'SL and 6'SL in the millimolar affinity range, 0.57 mM and 1.70 mM,

376    respectively. This is comparable to the affinity recently measured for the isolated *S.*

377    *pneumoniae* SpCBM40_NanC[37] against 3'SL (Kd ~ 1.5 mM) and 6'SL (Kd ~ 1.6 mM). Low

378    sialic acid affinity has also been proposed for *Cp*CBM40_NanJ from the *C. perfringens*

379    sialidase however this was not quantified[25]. Micromolar sialic acid affinity has been observed

380    for *C. perfringens* CpCBM40_NanI and *S. pneumoniae* SpCBM40_NanA[32,35]. Additional

381    electrostatic interactions with the sialic acid glycerol moiety may contribute to these unusual

382    affinities, in the case *Sp*CBM40_NanA via the introduction of a tryptophan in place of

383    *Rg*CBM40 Tyr210 (**Supplementary Fig. 3a, b**), and in the case of *Cp*CBM40_NanI via

384    Asn158, which approaches the binding site from a nearby loop extension (**Supplementary**

385    **Fig. 3c**). *Cp*CBM40_NanI also introduces additional water mediated interactions with the

386    galactose residues of bound 3'SL via a further loop extension (**Supplementary Fig. 1h**):

387    These are proposed to provide specificity for the corresponding sialic acid linkage[32]. A

388    corresponding extension is absent in *Rg*CBM40 leading to minimal observed interactions

389    between the protein and galactose (**Fig.1f, g, Supplementary Fig 1a**). Similar absence in

390    *Sp*CBM40_NanA suggests that these water-mediated interactions are not the defining

391    feature of high CBM40 sialic acid affinity.

392    Overall the binding epitopes of 3'SL and 6'SL, as determined by STD NMR, were in

393    agreement with the crystal structure, and confirmed the flexibility of the galactose and

394    glucose rings at the reducing end. Although the sialic acid moiety was the main recognition

395    element for the interaction with *Rg*CBM40, only weak binding was observed to Neu5Ac or

396    Neu5Gc monosaccharides. Sialic acid residues present in oligosaccharides are α-anomers.

397    However, in solution sialic acid adopts both α and β-anomeric configurations, as well as an

398    open chain conformation, with the β-anomer forming the dominant constituent[45]. In the

399    *Rg*CBM40 complex crystal structures, sialic acid is bound in the α-anomeric conformation,

400    allowing the axial C2 carboxylic acid moiety to form a conserved interaction with Arg204.

401    The *Rg*CBM40 preference for the minority α-anomer will incur a large entropic penalty. This

402    may provide a major contributory factor to the low observed monosaccharide affinity.

403    Thermodynamic analysis showed that the reaction is driven by enthalpy, with unfavorable

404    entropy (**Supplementary Table 2**), which is typical of interactions between CBMs and

405    saccharides[46].

406    The binding specificity of CBMs most commonly matches that of the appended catalytic

407    module[26,47]. We previously showed that the catalytic activity of *Rg*NanH is specific for α2-3-

408    linked sialic acid[30]. However, our glycan array and STD NMR data clearly showed that

409    *Rg*CBM40 can recognize a wide range of α2-3- and α2-6-sialic acid-linked oligosaccharides

410    which are commonly found in human GI mucins[12,21,23], suggesting an additional function.

411    More than 100 complex oligosaccharides were identified in mucins from human colonic

412    biopsies where most were mono-, di- or trisialylated[23]. *Rg*CBM40 bound Neu5Acα2-6Tn and

413    Neu5,9Ac$_2$α2-6Tn, Neu5Acα2-3TF and Neu5,9Ac$_2$-TF9Ac$_2$α2-3TF *but* not to the non-

414    sialylated forms; it also recognises Neu5Ac and acetylated Neu5Ac-linked Lac with α2-3 and

415    α2-6 linkage but shows a strict preference for Neu5Ac-linked LacNAc with α2-3 linkage, in

416    line with the increased expression of group Sd(a)/Cad related epitopes GalNAcβ1-

417    4(NeuAcα2-3)Gal along the length of the colon[12]. Despite the large diversity of structures,

418    the sigmoid MUC2 O-glycan repertoire and relative amounts in normal individuals is

419    relatively constant[23], suggesting their role in selecting a specific mucus-associated

420    microbiota.  Many bacterial species bind host tissues through protein-carbohydrate

421    interactions *via* a variety of cell-surface proteins and appendages. Although a wide number

422    of microbial lectins have been functionally and structurally characterized to date, especially

423    from pathogens, only a few carbohydrate-binding proteins present in gut bacteria which

424    interact with mucus have been structurally characterized[13,15]. Interactions between bacterial

425    adhesins from gut commensals and mucin glycans are generally of low affinity, in line with

426    the localization of these bacteria within the outer mucus layer[48,49]. Here we showed that

427    *Rg*CBM40 could recognize mucins with binding affinity increasing with sialic acid level.

428    Binding was highest towards human colonic MUC2, consistent with the increasing sialic acid

429    gradient along the GI tract from the small intestine to the colon in humans[21]. This study

430    demonstrates CBM40 mediating interaction to mucus, therefore expanding the repertoire of

431    bacterial adhesins to mucus. In addition to variations along the length of the GI tract, mucin

432    sialylation varies significantly between species, and thus could influence host species and

433    niche specificity of the gut symbionts. Interestingly, *Rg*CBM40 also showed binding to

434    Neu5Gc-containing oligosaccharides, albeit to lower affinity as compared to Neu5Ac-

435    oligosaccharides. Humans express predominantly Neu5Ac whereas Neu5Gc is expressed in

436     many non-human mammals[50]. Therefore, the ability of CBM from human gut commensal

437     bacteria to bind to Neu5Gc was unexpected. However, it cannot be excluded that *Rg*CBM40

438     mediates binding to dietary Neu5Gc-containing glycoproteins[51].

439     CBMs typically function to maintain carbohydrate-active enzymes (CAZymes) in proximity of

440     the substrate, thereby enhancing catalytic activity[26,46,52,53]. It has recently been suggested

441     that CBMs may play an additional role in the host–bacterium interaction by not only

442     mediating the attachment of CAZymes to glycans present on host tissues but by aiding the

443     adherence of the entire bacterium[27]. This would be particularly relevant to bacteria of the

444     human gut microbiota which are characterized by their large and diverse repertoires of CBM-

445     containing CAZymes[54]. Many CAZymes are known, or postulated to be, attached to the

446     bacterial cell surface[4]. Here, immunogold labeling confirmed the presence of *Rg*NanH on *R.*

447     *gnavus* ATCC 29149 cell-surface but not on *R. gnavus* E1. In addition, we showed that the

448     binding of *R. gnavus* ATCC 29149 to intestinal mucus was sialic acid mediated. The

449     potential avidity effect of CBM40-mediated binding of sialylated mucins *in vivo* (when

450     naturally present on the bacterial cell surface), may favor a mechanism by which CBM40

451     helps targeting the bacteria towards sialic acid rich regions of the GI tract, therefore

452     promoting bacterial colonization within the outer mucus layer. Our bioinformatics analyses of

453     bacterial genomes showed that *Rg*CBM40 canonical type domains are widespread among

454     Firmicutes, also reflecting the strong difference in CAZyme content and diversity between

455     the Firmicutes and Bacteroidetes phyla[54]. We thus propose a new role of CBMs in assisting

456     the tropism and spatial distribution of symbiotic bacteria among physical niches in the gut.

457

458

**Methods**

**Materials**

General chemicals including Neu5Ac were from Sigma (St Louis/MOI, US). Neu5Gcα2-3Lac Neu5Gcα2-6Lac, Neu5Ac-STn, Neu5Gc-STn and STFαOC3H6N3) were synthesised following published methodology[38,55]. Neu5Gc, 3'SL, 6'SL, 3'SGal, 6'SGal, 3'SLN, 6'SLN, were from Carbosynth. 2,7-anhydro-Neu5Ac was synthesised as previously reported[31]. Sialidase from *Clostridium perfringens* and *Salmonella typhimurium* LT2 were from New England Biolabs (Ipswich, MA US). Sialidase 0625 from *Akkermansia muciniphila* was a gift from WM de Vos[30]. Polyclonal antiserum against IMAC-purified His$_6$-*Rg*NanH[30] was raised in rabbits by BioGenes GmbH (Berlin, Germany) and provided at a titre of >1:200 000. Protease inhibitors benzamidine, N-ethylmaleimide, PMSF, sodium azide and soy bean inhibitor were from Sigma. Fluorescein labelled *Sambucus nigra* lectin (SNA-FITC) biotinylated SNA (SNA-biotin) and Vectashield were from Vector laboratories (Peterborough, UK). Streptavidin Alexa Fluor 488 conjugate was Thermo Fischer Scientific (Eugene/OR, US). Deuterium oxide (99.9% 2H) and Tris(hydroxymethyl-d3)amino-d2-methane (Tris-d11, 98% 2H) were from Sigma. Mouse monoclonal anti-His-HiLyte Flour 555 antibody was obtained from LifeSpan BioSciences (Seattle/WA, US). Blocking reagent was from Perkin Elmer (Boston/MA, US). Rabbit Mucin 2 antibody H-300 was from Santa Cruz (Dallas/TX, US, SC-15334), Goat anti-Rabbit IgG Secondary Antibody, Alexa Fluor 488 (A11034) and Goat anti-Rabbit IgG Secondary Antibody Alexa Fluor 594 (A11037) from Thermo Fischer Scientific. DAPI was from Life Technologies, O.C.T. Compound from VWR and Hydromount from National Diagnostics (Atlanta/GA, USA).

**Expression and purification of *Rg*CBM40 and *Rg*NanH**

Using the full-length sequence encoding *Rg*NanH in pOPINF from *R. gnavus* strain ATCC 29149 as a template[30], *Rg*CBM40 (residues 50–237), *Rg*NanH (residues 26–723) and *Rg*GH33 (residues 243–723) were cloned into the pEHISTEV vector[56] using the primers listed in **Supplementary Table 4.** Protein expression and purification of *Rg*CBM40 and *Rg*NanH was similar to that of *Rg*GH33[30]. Points of divergence are indicated below. For protein expression, recombinant plasmids were transformed into *E. coli* BL21 Rosetta (DE3) (Novagen, NJ, US). A single colony was used to inoculate a 10 ml Luria Bertani (LB) medium pre-culture, which was incubated overnight under shaking at 200 rpm (at 30 °C for crystallisation and protein size determination or at 37 °C for all other protein assays). The pre-culture was used to inoculate 500 ml of auto induction medium (Formedium, Norfolk, UK), which was incubated under shaking at 37 °C for 3 h followed by 60 h incubation at 16 °C. All cultures were inoculated with 50 μg ml$^{-1}$ kanamycin.

15

496    For crystallisation and protein size determination, cells were harvested by centrifugation,

497    resuspended in phosphate buffered saline (PBS, 150 mM sodium chloride, 10 mM sodium

498    phosphate, pH 7.4) for *Rg*CBM40 and in 20 mM Tris-HCl pH 7.5, 50 mM NaCl for *Rg*NanH,

499    supplemented with DNase I (20 µg ml$^{-1}$) and cOmplete protease inhibitor mixture tablets

500    (Roche, Welwyn Garden City, UK), and lysed using a constant flow cell disrupter. Insoluble

501    components were removed by centrifugation and filtration through a 0.22 µm pore size

502    syringe driven filter (Millipore, NJ, US). Soluble lysate was loaded onto a nickel-Sepharose

503    column (GE Healthcare, Little Chalfont, UK) overnight at 4 °C. The sample was then washed

504    extensively with lysis buffer supplemented with 5 mM imidazole for *Rg*CBM40 and with

505    150 mM imidazole for *Rg*NanH and was eluted using lysis buffer supplemented with 50 mM

506    imidazole for *Rg*CBM40 and with 300 mM imidazole for *Rg*NanH. The sample was then

507    dialysed into lysis buffer and cleaved of its six-histidine tag using Tobacco Etch Protease at

508    a mass ratio of 1:50 overnight at 4 °C. Finally, the gel filtration step using a Sephacryl S-100

509    column (GE Healthcare) was performed using 20 Tris pH 7.5 with 50 mM NaCl. The purified

510    *Rg*CBM40 was crystallised as described below. To determine the size in solution of

511    *Rg*NanH, size exclusion chromatography with multi angle light scattering (SEC-MALS) was

512    performed using an NGC chromatography system (Biorad, Hercules, CA,US) equipped with

513    a DAWN HELEOS II MALS detector (Wyatt technology, Haverhill, UK) and an Optilab T-rEX

514    differential Refractive Index detector (Wyatt Technology). The data were analysed using

515    ASTRA (Wyatt Technology).

516    For all other protein assays, the cell pellets were resuspended in Bug buster-HT (Merck,

517    Kenilworth, NJ, US) with the supplied lysozyme and lysed by shaking in this solution for 1 h

518    at room temperature. Insoluble material was removed by centrifugation at 4 °C, 3320 *g* for

519    25 min and the supernatant was dialysed into desalting buffer (50 mM Tris-HCl, 150 mM

520    NaCl, pH 7.8 containing 10 mM imidazole for *Rg*GH33 and *Rg*NanH and no imidazole for

521    *Rg*CBM40, the difference is due to the poor binding of the His$_6$-tag of *Rg*CBM40 to the nickel

522    column) to remove the Bug buster-HT. Again insoluble material was removed by

523    centrifugation as above, except at 8 000 *g*. Purification of the soluble lysate was loaded onto

524    the immobilized metal ion affinity chromatography (IMAC column, His-bind, Novagen) in

525    binding buffer (desalting buffer with the addition of 10 mM imidazole) using the Akta Express

526    (GE Healthcare). The protein was eluted with binding buffer containing 500 mM imidazole

527    and then immediately desalted into desalting buffer. The partially purified protein was

528    concentrated using 3.5 kDa MWCO spin columns (Sartorius, Gottingen, Germany) prior to

529    gel filtration again with the Akta Express in desalting buffer (see above) on a Superdex 75

530    column (GE Healthcare). Purity of the proteins was assessed throughout by SDS-PAGE

531    using the Novex system (Thermo Fisher Scientific).

532

**Site-directed mutagenesis**

533 

534 Site directed mutagenesis of *Rg*GH33 to introduce the D282A mutation in the active site was

535 carried out using the QuikChange kit, following the manufacturer's instructions, Agilent

536 (Santa Clara, CA, US). Site-directed mutants of *Rg*CBM40; I95A, Y116A, E126A, R128A,

537 R204A and double mutant R128A/R204A, were obtained from NZyTech (Lisbon, Portugal).

538 The primers are listed in **Supplementary Table 4**. The integrity of the *Rg*GH33 and

539 *Rg*CBM40 mutants was checked by circular dichroism (CD).

540 

**Circular dichroism**

541 

542 CD spectra were recorded using a JASCo J-700 spectropolarimeter, under the following

543 conditions: 20 nm/min scan speed, bandwidth 1 nm, response 2 s, 5 points/nm and 4

544 accumulations. Far-UV spectra (260-180nm) were recorded in a 0.1 mm pathlength cell. The

545 spectropolarimeter was calibrated using camphorsulphonic acid (Sigma). The protein was

546 extensively dialysed into 10 mM sodium phosphate buffer, pH 6.5 and a buffer only control

547 was subtracted from all spectra using the molar CD factor calculated as follows: (113 x 30 x

548 $10^{-6}$)/ [conc(mg ml$^{-1}$) x pathlength (cm)].

549 

**Protein crystallisation**

550 

551 The final crystallisation condition was 0.2 M ammonium chloride with 20% PEG 8000. The

552 drop contained 0.5 µl protein solution at 25 mg ml$^{-1}$ and 0.5 µl reservoir solution, initial

553 crystals grew in four weeks and growth time was improved significantly using micro

554 seeding[57]. Crystals were cryoprotected using the crystallisation condition supplemented with

555 25% (w/w) glycerol. To achieve crystal structures in complex with 3'SL and 6'SL the crystals

556 were grown in crystallisation condition supplemented with 20 mM ligand followed by a 60

557 min soak in crystallisation condition supplemented with 100 mM ligand immediately prior to

558 cryoprotection and mounting.

559 

**Solving the crystal structure**

560 

561 X-ray diffraction experiments were performed at 100 K. Data were collected using a Rigaku

562 MSC Micromax 007 HF X-ray source, with a fixed wavelength of 1.542 Å, and a Saturn 944+

563 CCD detector. Sweeps were indexed and integrated separately and then scaled together

564 within the HKL2000 data processing package[58]. Phasing was performed by Phaser[59] within

565 the CCP4 package[60] using the CBM40 of the *M. decora* sialidase NanL (*Md*CBM40_NanL)

566 (PDB 2SLI)[33] as the molecular replacement model. The model was refined using iterative

567 cycles of Refmac5[61] and Coot[62]. The PDB REDO server was used to optimize the

568 refinement parameters[63]. The model was validated using the Molprobity server[64]. Paired

569     refinement performed by the PDB REDO server indicated that the models were improved by

570     the inclusion of high resolution, low completeness data for the 3'SL and 6'SL complexes[65].

571     For an illustrative stereo image of a portion of the electron density map, see **Supplementary**

572     **Fig. 13**.

573

574     **Isothermal titration calorimetry**

575     ITC experiments were performed using the PEAQ-ITC system (Malvern, Malvern, UK) with a

576     cell volume of 200 µl. Prior to titration protein samples were exhaustively dialyzed into PBS.

577     The ligand was dissolved in the dialysis buffer. The cell protein concentration was 115 µM

578     (except for mutant I95A where it was 173 µM and the wild type interaction with 6'SL where it

579     was 230 µM) and the syringe ligand concentration was 10 mM (25 mM for Neu5Ac).

580     Controls with titrant (sugar) injected into buffer only were subtracted from the data.  Analysis

581     was performed using Malvern software, using a single binding site model. The stoichiometry

582     of binding sites was set to 1.0 as this was evident from the crystal structure. Quantitative and

583     most qualitative experiments were carried out in triplicate.

584

585     **STD NMR experiments**

586     $^1$H and $^{13}$C resonance assignment for all the sugars was performed on the bases of 1D $^1$H,

587     2D DQF-COSY, TOCSY, HSQC and NOESY experiments run on the free ligands in

588     unbuffered $D_2O$, pH 7.0. For STD NMR experiments, all the samples consisted of 1 mM

589     sialoglycans and 50 µM *Rg*CBM40 (WT or I95A mutant) in $D_2O$ buffer solution of 10 mM

590     Tris-$d_{11}$ pH 7.8 and 100 mM NaCl (ligand : protein ratio 20 : 1). An STD pulse sequence that

591     included 2.5 ms and 5 ms trim pulses and a 3 ms spoil gradient was used. Saturation was

592     achieved applying a train of 50 ms Gaussian pulses (0.40 mW) on the f2 channel, at 0.60

593     ppm (on-resonance experiments) and 40 ppm (off-resonance experiments). The broad

594     protein signals were removed using a 40 ms spinlock (T1ρ) filter. All the experiments were

595     recorded at $^1$H frequency of 800.23 MHz on a Bruker Avance III spectrometer equipped with

596     a 5 mm probe TXI 800 MHz H-C/N-D-05 Z BTO, at 288 K. For all the sialoglycans in the

597     presence of *Rg*CBM40, an STD experiment with a saturation time of 2 s and a relaxation

598     delay of 5 s was performed, as a first test for binding. For the confirmed binders, the STD

599     NMR experiments were carried out at different saturation times (0.5, 1, 2, 3, 4 and 5 s) with

600     1K scans and relaxation delay of 5 s, in order to obtain the binding epitope mapping. The

601     resulting build-up curves for each proton were fitted mathematically to a mono-exponential

602     equation ($y=a*[1-exp(b*x)]$), from which the initial slopes ($a*b$) were obtained. For each

603     ligand, the binding epitope mapping was obtained by dividing the initial slopes by the one of

604     the H7 proton of the corresponding sialic acid ring, to which an arbitrary value of 100% was

605     assigned. This normalization of the STD values allows the comparison across all the

606   sialoglycans.

607

**Structure-based sequence alignment and bioinformatics analyses**

609   A structural alignment of *Rg*CBM40 was carried out with all CBM40 structures available to
610   date (see Results and **Supplementary Methods**). This served as a basis for producing an
611   alignment including both canonical and *Vibrio* type CBM40 sequences to create a profile
612   Hidden Markov Model (pHMM) using the HMMER3 software (http://hmmer.org/)
613   (**Supplementary Fig. 14**), intended to detect both types simultaneously and ensure that hit
614   sequences of both types are thus properly aligned for subsequent comparative analysis.
615   Additionally, we created pHMMs corresponding to the canonical-only and *Vibrio* type-only
616   CBM40 sequences of this alignment, to resolve the type of each hit. Protein domain
617   databases such as Pfam[66] currently characterize the canonical CBM40 as a sequence family
618   belonging to a larger superfamily ("clan"), and some individual domains make good matches
619   to more than one related family, i.e. including non-CBM40 such as "Concanavalin A-like
620   lectin/glucanases" (in contrast, no Pfam domain clearly defines the *Vibrio* CBM40). We
621   therefore also used the corresponding Pfam pHMMs, as well as our own, to search all
622   available (177 million) protein sequences from annotated NCBI prokaryote genomes, using
623   HMMER3. Where individual hit domains matched multiple pHMMs, we compared scores to
624   identify and discard hits which might be better regarded as related, non-CBM40 domains.
625   The remaining CBM40 proteins were screened for the presence of the sialidase domain
626   (GH33) and IT-sialidase, as previously described[30]. We reduced this to a nonredundant set
627   (**Supplementary Methods**) for further analysis. A detailed phylogenetic analysis is beyond
628   the scope of this study, but we estimated evolutionary distances between these 51
629   representative sequences using fprotdist in EMBASSY-PHYLIP[67,68] from which the tree was
630   calculated by neighbour-joining (fneighbor). All sites were included in the analysis, using the
631   PMB model with a uniform rate of evolution. This was repeated on 1,000 replicate datasets
632   produced by bootstrap resampling (fseqboot; consensus tree produced by fconsense). The
633   figure was produced with FigTree (http://tree.bio.ed.ac.uk/software/figtree/). Bioinformatics
634   analyses were performed using the Gut Health and Food Safety Linux servers at Quadram
635   Institute Bioscience.

636

**Glycan microarray screening**

638   Glycan microarrays were fabricated using epoxide-derivatized slides as previously described
639   (38). Printed glycan microarray slides were blocked by ethanolamine, washed and dried.
640   Slides were then fitted in a multi-well microarray hybridization cassette (AHC4X8S, ArrayIt,
641   Sunnyvale, CA, USA) to divide into 8 subarrays. The subarrays were blocked with ovalbumin
642   (1% w/v) in PBS (pH 7.4) for 1 h at room temperature, with gentle shaking. Subsequently,

19

643    the blocking solution was removed and diluted protein samples of *Rg*CBM40 and *Rg*GH33

644    D282A with various concentrations were added to each subarray. After incubating the

645    samples for 2 h at room temperature with gentle shaking, the slides were washed. Diluted

646    anti-His-HiLyte Flour 555 antibodies in PBS were added to the subarrays, incubated for 1 h

647    at room temperature, washed and dried. The microarray slides were scanned by Genepix

648    4000B microarray scanner (Molecular Devices Corp., Union City, CA, USA). Data analysis

649    was performed using Genepix Pro 7.0 analysis software (Molecular Devices Corp.). It is

650    important to note that glycans on the array with sialic acid O-acetyl groups undergo gradual

651    losses of these labile ester groups. Therefore, definitive conclusions about 9-O-acetylation

652    are only possible in instances wherein binding is exclusively to the O-acetylated sialoglycan

653    spot, and not to the corresponding non-O-acetylated spot.

654

655    **RgCBM40 binding to mucus-producing cells**

656    The binding of *Rg*CBM40 to mucus-producing LS174T cell line (80% confluent, passage 12)

657    was performed by incubating the cells with 150 µg ml$^{-1}$ *Rg*CBM40 in cell culture medium for

658    2 h at 37 °C. Control samples were incubated with cell culture medium only. The cells were

659    then washed with PBS, fixed in methacarn (60% dry methanol, 30% chloroform and 10%

660    acetic acid) and washed in PBS containing 0.05% bovine serum albumin (BSA).  Blocking

661    was done with TNB buffer (0.5% w/v blocking reagent in 100 mM Tris-HCl pH 7.5, 150 mM

662    NaCl) supplemented with 5% goat serum. The *Rg*CBM40 binding was detected with custom-

663    made rabbit *Rg*NanH antiserum diluted 1:100 in PBS and goat anti-rabbit antibody diluted

664    1:400 in PBS. The same antibodies were used for negative control sample (*Rg*CBM40-free).

665    In the lectin control sample, SNA-biotin (incubated at 75 µg ml$^{-1}$) was detected with

666    streptavidin conjugate (2.5 µg ml$^{-1}$). MUC2 was detected with rabbit Mucin 2 antibody diluted

667    1:50 in PBS and goat anti-rabbit antibody diluted 1:200 in PBS. The cells were

668    counterstained with DAPI and mounted in Vectashield. The slides were imaged using a

669    Zeiss Axio Imager 2 microscope.

670

671    **RgCBM40 and *R. gnavus* binding to intestinal tissue**

672    To assess the binding of *Rg*CBM40 to intestinal tissue sections, colon of wild type C57BL/6

673    mouse was washed with PBS, fixed in methacarn, embedded in O.C.T. compound and cut

674    into 8 µm sections. Access to mouse tissues was carried out under the Animal Welfare and

675    Ethical Review Body of University of East Anglia's establishment licence (according to Home

676    Office requirements).  Tissue sections were washed in PBS containing 0.05% BSA and

677    blocked with TNB buffer (0.5% w/v blocking reagent in 100 mM Tris-HCl pH 7.5, 150 mM

678    NaCl) supplemented with 5% goat serum. The slides were then washed in PBS 0.05% BSA,

679    followed by 2 h incubation of 150 µg ml$^{-1}$ *Rg*CBM40 in PBS at 37 °C. Control tissue sections

680  were incubated in PBS only. After washes in PBS with 0.05% BSA, the binding of *Rg*CBM40

681  was detected with custom-made rabbit *Rg*NanH antiserum (diluted 1:100 in TNB buffer) and

682  goat anti-rabbit antibodies (diluted 1:200 in PBS). Negative control sample (*Rg*CBM40-free)

683  was also incubated with these primary and secondary antibodies. Muc2 was detected with

684  Mucin 2 antibody diluted 1:100 in TNB buffer and goat-anti rabbit antibody diluted 1:200 in

685  PBS. In lectin controls SNA-FITC was incubated at 4 µg ml$^{-1}$. The sections were

686  counterstained with DAPI and mounted in Hydromount mounting medium. The slides were

687  imaged using Zeiss an Axio Imager 2 microscope. To assess the binding specificity of

688  *Rg*CBM40 to sialylated structures, the tissue sections were pre-treated with sialidase.

689  Briefly, saponification was performed to make the enzymatic digestion of mouse colonic

690  tissue sections effective[69]. The sections were treated with 0.5% KOH in 70% ethanol for 15

691  min at room temperature. After three PBS washes, 500 U ml$^{-1}$ sialidase from *Clostridium*

692  *perfringens* in GlycoBuffer 1 (New England Biolabs) was added and incubated for 14 h at 37

693  °C. Sections were incubated in sialidase-free GlycoBuffer 1 under the same experimental

694  conditions and used as a control of sialidase digestion to assess the binding *Rg*CBM40 and

695  SNA to tissue sections as described above.

696  To assess the binding of *R. gnavus* to intestinal tissue sections, colon of wild type C57BL/6

697  mouse was washed with PBS, fixed in methacarn, embedded in O.C.T. compound and cut

698  into 12 µm sections. Tissue sections were washed in PBS, then incubated with SNA in PBS

699  at 20 µg ml$^{-1}$ for 1 h. Prior to incubation with bacteria, the slides were washed with PBS. *R.*

700  *gnavus* ATCC 29149 was cultured anaerobically in BHI-YH media for 24 h as previously

701  described[29]. The culture was then then used to inoculate YCFA media supplemented with

702  3'SL at a concentration of 7 mg ml$^{-1}$, and cultured for 20 h. The bacteria were then washed

703  twice with fresh YCFA, and resuspended at an OD of 1. The tissue sections were then

704  transferred in a humid chamber to the anaerobic cabinet, and the bacteria incubated on the

705  sections for 1 h at 37$^{o}$C. The slides were then washed twice with YCFA and fixed with 4%

706  paraformaldehyde in PBS for 15 min. The slides were transferred out of the anaerobic

707  cabinet, then washed with PBS and blocked with TNB buffer (0.5% w/v blocking reagent in

708  100 mM Tris-HCl pH 7.5, 150 mM NaCl) supplemented with 5% goat serum. The presence

709  of *R. gnavus* and Muc2 was detected with custom-made rabbit *Rg*NanH antiserum (diluted

710  1:100) and Mucin 2 antibody (1:100), respectively. Goat anti-rabbit antibodies (diluted 1:500)

711  were used for immunodetection. The sections were counterstained with DAPI and mounted

712  in Prolong gold anti-fade mounting medium. The slides were imaged using Zeiss an Axio

713  Imager 2 microscope, using a x63 objective.

714

715  **Mucin purification**

716 Culture media from LS174T cell line were freeze-dried before extraction of MUC2. After

717 freeze-drying, samples were solubilised overnight in 6 M guanidine chloride (GuCl) buffer

718 containing protease inhibitors (7.95 mM EDTA, 12.25 mM benzamidine, 6.25 mM *N*-

719 ethylmaleimide, 1.25 mM PMSF, 3.75 mM sodium azide, 0.1 mg/ml soy bean inhibitor).

720 Samples were centrifuged at 18 500 *g*. The pellet was reduced with dithiothreitol (DTT) at 10

721 mM for 4 h at 45 °C and alkylated with 25 mM iodoacetamide overnight before dialysis

722 against 50 mM ammonium bicarbonate. The same protocol was followed for purifying

723 mucins from the scraped mucus from small intestine and colon of mouse models. The

724 supernatants containing soluble mucins were diluted in 4 M guanidinium chloride (GuCl) with

725 phosphate buffered saline (PBS) and adjusted with cesium chloride at 1.4 g ml$^{-1}$ density.

726 Supernatants were subjected to an ultracentrifugation (Beckman, Brea, US) at 234 000 *g* for

727 72 h at 20 °C. Fractions of 1 ml were collected and weighed. Fractions between 1.35 and

728 1.45 g ml$^{-1}$ were kept and dialysed against 50 mM ammonium bicarbonate. These fractions

729 contained the purified mucins.

730

731 **Release of oligosaccharides from mucin**

732 The mucins were subjected to β-elimination under reductive conditions (0.1 M sodium

733 hydroxide, 1 M sodium borohydride) for 20 h at 45 °C. The reaction was stopped by adding

734 Dowex 50 x 8 (Sigma) and filtered before being co-evaporated with methanol 3 times.

735 Remaining salts were removed by Carbograph (Grace, Columbia, US).

736

737 **Permethylation of O-glycans**

738 Permethylation was performed on released *O*-glycans from the different mucins samples.

739 Samples were solubilized in 200 µl dimethyl sulfoxide. Then sodium hydroxide (trace of

740 powder) and 300 µl iodomethane were added in anhydrous conditions and the samples

741 vigorously shaken at room temperature for 90 min. The permethylation reaction was stopped

742 by addition of 1 ml acetic acid (5% vol/vol). Permethylated *O*-glycans were purified on a

743 Hydrophilic-Lipophilic Balanced (HLB) Oasis cartridge (Waters, Milford, US). Briefly,

744 cartridges were activated by methanol, equilibrated with methanol:water (5:95, vol:vol), and

745 samples loaded onto the cartridges. Cartridges were washed by methanol:water (5:95,

746 vol:vol) and the permethylated *O*-glycans eluted by methanol.

747

748 **Analysis of permethylated O-glycans by mass spectrometry**

749 MALDI-TOF and TOF/TOF-MS data were acquired using the Bruker Autoflex analyzer mass

750 spectrometer (Applied Biosystems, Foster City, CA, US) in the positive-ion and reflectron

751 mode by using 2,5-dihydroxibenzoic acid (DHB; Sigma; 10 mg ml$^{-1}$ in 70:30 methanol:water)

752 as the matrix. The relative quantification of sialylation on mucins was calculated based on

753    the sum of all areas of mass peaks corresponding to sialylated structures divided by the sum

754    of all areas of mass peaks corresponding to defined O-glycans.

755

756    **Enzyme Linked Immunosorbent Assay**

757    *Rg*CBM40 binding to purified mucins was tested by ELISA.  Mucins (100 μl of 10 μg ml$^{-1}$)

758    were immobilised onto a high binding 96 well plate (Greiner, Stonehouse, UK) overnight at 4

759    °C. All subsequent steps were carried out for 1 h at room temperature. The plates were

760    blocked with 3% (w/v) BSA, incubated with *Rg*CBM40 (500 μg ml$^{-1}$), followed by an

761    incubation with 1:5 000 anti-*Rg*NanH (raised in rabbit, Biogenes) then with 1:5 000 anti-

762    rabbit secondary antibody (raised in donkey) conjugated to peroxidase (GE Healthcare).

763    Between each step the plate was washed with 3 x 300 ul of PBS containing 0.05% (v/v)

764    Tween 20 (PBST). Prior to detection, an additional wash step and 30 sec incubation with

765    PBST was carried out. Binding was detected using tetramethylbenzidine (TMB) visualisation

766    solution (Biolegend, San Diego, CA, US) which was incubated for 15 min. The reaction was

767    stopped by addition of 2 M $H_2SO_4$ and absorbance measured at 450 nm using a plate-reader

768    (Bench Marl Plus, Biorad), subtracting background readings at 570 nm. Negative controls

769    including no *Rg*CBM40 (subtracted from $A_{450}$ value), no primary or no secondary antibody

770    were carried out in parallel. For comparison between plates, values were normalised to the

771    reading for LS174T MUC2 which was arbitrarily set at 100%. For enzymatic treatment of the

772    mucin, LS174T MUC2 (2 mg ml$^{-1}$) was incubated with sialidases (2 μg ml$^{-1}$) overnight at 4 °C

773    on a rotary wheel prior to immobilization on the plate. For chemical treatment of mucin,

774    LS174T MUC2 was incubated with 0.1 M trifluoroacetic acid (TFA) at 80 °C for 1 h, dialysed

775    against ammonium bicarbonate (50 mM), lyophilised and redissolved in $H_2O$. For the

776    competition assays, *Rg*CBM40 was incubated with 1 mM of free sugar overnight at 4°C on a

777    rotary wheel prior to addition to the ELISA plate as above. Experiments were carried out in

778    triplicate.

779

780    **HPAEC-PAD analyses**

781    The substrates, 3'SL (500 μM, 8.5 nM enzyme), 3'SLX (Neu5Ac form), 500 μM, 80 nM

782    enzyme) or LS174T MUC2 (0.9 mg ml$^{-1}$, 1.5 nM enzyme) were incubated with *Rg*NanH or

783    *Rg*GH33 at 37 °C in 20 mM sodium phosphate buffer, pH 6.5. BSA (0.1 mg ml$^{-1}$) was

784    included in the oligosaccharide reactions. Control reactions without enzyme were also

785    carried out in parallel. Aliquots of reaction were removed and the reaction terminated by

786    boiling for 20 min. For LS174T MUC2, the released sugars were removed using 5 kDa

787    MWCO spin columns and the remaining mucin subjected to acid hydrolysis; the samples

788    were incubated with 0.1 M HCl at 80 °C for 1 h, dried under vacuum and resuspended in

789     $H_2O$ at 1 mg ml$^{-1}$. The amount of Neu5Ac remaining on the mucin was quantified by

790     comparing the peak size for Neu5Ac with an internal standard of 2-keto-3-deoxynononic acid

791     (Kdn). The reaction products for all substrates were filtered with 0.22 μm spin tubes prior to

792     analysis by HPAEC-PAD (Dionex ICS-5000, Thermo Fisher Scientific). An internal standard

793     of fucose (50 μM) was used for 3'SL and 3'SLX. For 3'SL, a Carbo-Pac PA1 column

794     (Thermo Fisher Scientific) was used with a 6 min isocratic gradient of 100 mM sodium

795     hydroxide, 100 mM sodium acetate followed by a 10 min washing step with 100 mM sodium

796     hydroxide, 200 mM sodium acetate and 10 min re-equilibration with 100 mM sodium

797     hydroxide, 100 mM sodium acetate. For 3'SLX, a Carbo-Pac PA100 was used with 5 min at

798     100 mM sodium hydroxide, a gradient of 0–50 mM sodium acetate over 5 min, followed by a

799     gradient of 50–225 mM sodium acetate. The column was then cleaned with 500 mM sodium

800     acetate for 5 min and re-equilibrated for 15 min at 100 mM sodium hydroxide. For analysis of

801     the acid hydrolysis products of MUC2, a Carbo-Pac PA10 was used with a gradient of 70–

802     300 mM sodium acetate with 100 mM sodium hydroxide over 10 min, a brief (1 min) period

803     of 300 mM sodium acetate followed by a decrease (over 1 min) to 70 mM sodium acetate

804     and 15 min re-equilibration at 70 mM sodium acetate. All columns were protected with their

805     respective guard columns, except for the mucin analysis where an amino-guard column was

806     used.

807

808     **Western blotting**

809     *R. gnavus* strains were grown to stationary phase and cells pelleted by centrifugation for 10

810     min at 3 000 *g* at 4 °C. The supernatant was collected and the extracellular proteins

811     concentrated 50-fold using a 10-kDa MWCO Amicon Ultra-0.5 Centrifugal Filter (Millipore,

812     Watford, UK). The cell pellet was re-suspended in 20 μl PBS with an equal bead (100 μm

813     glass beads) volume added and samples vortexed at full speed three times for 2 min with 2

814     min rest intervals on ice. The volume was made up to 17 μl per mg wet cell weight with PBS

815     and vortexed at full speed again for 2 min. The beads were removed by allowing them to

816     settle under gravity and the remaining samples centrifuged for 30 min at 17 000 *g* at 4 °C.

817     The supernatant containing the soluble cytosolic proteins was collected and concentrated

818     10-fold using a 10-kDa MWCO Amicon Ultra-0.5 Centrifugal Filter. The remaining pellet was

819     dissolved in 1.7 μl per mg wet cell weight digestion buffer (50 mM Tris-HCl (pH 8.0), 5 mM

820     $MgCl_2$, 5 mM $CaCl_2$, 10 mg ml$^{-1}$ Hen Egg White Lysozyme (Sigma), and incubated at 37°C

821     for 3 h. The samples were centrifuged for 30 min at 17 000 *g* at 4 °C, and the supernatant

822     containing the cell wall associated proteins collected. Samples were analysed on duplicate

823     NuPAGE Novex 4–12% Bis-Tris gels, one gel was stained with InstantBlue stain (Expedeon,

824     Swavesey, UK) and the other gel blotted onto a PVDF membrane using X-cell II Blot module

825    (Thermo Fisher Scientific), according to manufacturer's instructions. Membranes were

826    blocked with 3% BSA in PBST for 3 h, and then incubated with the custom-made anti-

827    *Rg*NanH antibody raised in rabbit diluted 1:5000 in 1% BSA in PBST overnight. Blots were

828    washed in PBST, then incubated with anti-rabbit IgG antibody (Sigma) diluted 1:7 500 in 1%

829    BSA in PBST for 2 h. After washing three times in PBST, the blots were incubated using a

830    visualisation solution (10 ml of 0.1 M Tris-HCl (pH 9.6), 40 µl of 1 M $MgCl_2$, 20 µl of nitroblue

831    tetrazolium, and 10 µl of 5-Bromo-4-Chloro-3-Indolyl phosphate, Sigma) for up to 15 min,

832    and washed in distilled water to stop the development of the signal.

833

834    **Immunogold labelling of whole bacterial cells**

835    *R. gnavus* strains were grown to stationary phase and cells pelleted by centrifugation for 10

836    min at 3 000 *g* at 4 °C before being resuspended in PBS. A small drop of concentrated *R.*

837    *gnavus* cell suspension was applied to a formvar/carbon coated gold TEM grid (Agar

838    Scientific, Stansted, UK) and left for 1 min. The bacteria on the grids were vapour fixed by

839    placing the grids in a sealed Petri dish with a small cap-full of 25% glutaraldehyde (Agar

840    Scientific) for 2 h. The grids were floated on drops of 50 mM Glycine/PBS for 15 min

841    followed by floating on drops of Aurion blocking buffer (Aurion, Wageningen, The

842    Netherlands) for 30 min. The grids were then washed five times for 5 min with 0.1% BSA-C

843    (Aurion) in PBS. Grids were incubated in anti-*Rg*NanH antibody raised in rabbit diluted

844    1:2000 with 0.1% BSA-C/PBS or in a control solution of 0.1% BSA-C/PBS overnight at 4°C.

845    The grids were washed five times for 5 min with 0.1% BSA-C/PBS. Grids were then

846    transferred to a 1/50 dilution of goat-anti-rabbit antibody conjugated with 10 nm gold balls

847    (Agar Scientific) in 0.1% BSA-C/PBS and incubated for 2 h at room temperature. The grids

848    were washed five times for 5 min with 0.1% BSA-C/PBS, followed by three 5 min washes in

849    PBS only. The grids were refixed by immersing them in 2% glutaraldehyde/PBS for 1.5 h

850    followed by three 5 min PBS washes and three 5 min distilled water washes before the grids

851    were carefully blotted and dried. The grids were examined and imaged in a FEI Tecnai G2

852    20 Twin transmission electron microscope at 200 kV.

853

854    **Statistical analysis**

855    One-way ANOVA model analyses were used to assess the binding of *Rg*CBM40 to purified

856    mucins by ELISA. When the effect of the factor was found to be significant (p value < 0.05)

857    and its number of levels greater than 2, a Tukey test was used to assess the significance of

858    the difference between multiple means. Statistical analyses were performed using the

859    software SAS 9.4 (NC, USA).

860

**Data availability**

Atomic coordinates have been deposited in the Protein Data Bank (www.rcsb.org) with

accession codes: unbound, 6ER2; 3'SL bound, 6ER3, 6'SL bound; 6ER4. All other relevant

data are available from the authors.


**References**

1. Sekirov, I., Russell, S. L., Antunes, L. C. M. & Finlay, B. B. Gut Microbiota in Health and Disease. *Physiol. Rev.* **90,** 859-904 (2010).

2. Donaldson, G. P., Lee, S. M., & Mazmanian, S. K. Gut biogeography of the bacterial microbiota. *Nat. Rev. Microbiol.* **14,** 20-32 (2016).

3. Martens, E. C., Chiang, H. C. & Gordon, J. I. Mucosal glycan foraging enhances fitness and transmission of a saccharolytic human gut bacterial symbiont. *Cell Host Microbe* **4,** 447-457(2008).

4. Flint, H. J., Scott, K. P., Duncan, S. H., Louis, P. & Forano, E. Microbial degradation of complex carbohydrates in the gut. *Gut Microbes* **3,** 289-306 (2012).

5. Johansson, M. E. V., Larsson, J. M. H. & Hansson, G. C. The two mucus layers of colon are organized by the MUC2 mucin, whereas the outer layer is a legislator of host–microbial interactions. *Proc. Natl. Acad. Sci. U. S. A.* **108,** 4659-4665 (2011).

6. McGuckin, M. A., Lindén, S. K., Sutton, P. & Florin, T. H. Mucin dynamics and enteric pathogens. *Nat. Rev. Microbiol.* **9,** 265-278 (2011).

7. Manichanh, C., Borruel, N., Casellas, F. & Guarner, F. The gut microbiota in IBD. *Nat. Rev. Gastroenterol. Hepatol.* **9,** 599-608 (2012).

8. Sheng, Y. H., Hasnain, S. Z., Florin, T. H. J. & McGuckin, M. A. Mucins in inflammatory bowel diseases and colorectal cancer. *J. Gastroenterol. Hepatol.* **27,** 28-38 (2012).

9. Li, H. *et al.* The outer mucus layer hosts a distinct intestinal microbial niche. *Nat. Commun.* **6,** 8292 (2015).

10. Ouwerkerk, J. P., de Vos, W. M. & Belzer, C. Glycobiome: Bacteria and mucus at the epithelial interface. *Best Pract. Res. Clin. Gastroenterol.* **27,** 25-38 (2013).

11. Jensen, P. H., Kolarich, D. & Packer, N. H. Mucin-type O-glycosylation--putting the pieces together. *FEBS J.* **277,** 81-94 (2010).

12. Robbe, C., Capon, C., Coddeville, B. & Michalski, J. C. Structural diversity and specific distribution of O-glycans in normal human mucins along the intestinal tract. *Biochem. J.* **384,** 307-316 (2004).

13. Juge, N. Microbial adhesins to gastrointestinal mucus. *Trends Microbiol.* **20**, 30-39 (2012).

14. Tailford, L.E., Crost, E.H., Kavanaugh, D. & Juge, N. Mucin glycan foraging in the human gut microbiome. *Front. Genet.* **6,** 81 (2015).

15. Etzold, S. & Juge, N. Structural insights into bacterial recognition of intestinal mucins. *Curr. Opin. Struct. Biol.* **28,** 23-31 (2014).

16. Ng, K. M. *et al.* Microbiota-liberated host sugars facilitate post-antibiotic expansion of enteric pathogens. *Nature* **502,** 96-99 (2013).

17. Tong, M. *et al.* Reprograming of gut microbiome energy metabolism by the FUT2 Crohn's disease risk polymorphism. *ISME J.* **8,** 2193-2206 (2014).

905   18. Bergstrom, K. S. & Xia, L. Mucin-type O-glycans and their roles in intestinal
906   homeostasis. *Glycobiology* **23,** 1026-1037 (2013).

907   19. Lewis, A. L. & Lewis, W. G. Host sialoglycans and bacterial sialidases: a mucosal
908   perspective. *Cell. Microbiol.* **14,** 1174-1182 (2012).

909   20. Juge, N., Tailford, L. & Owen, C. D. Sialidases from gut bacteria: a mini-review.
910   *Biochem. Soc. Trans.* **44,** 166-175 (2016).

911   21. Robbe, C. *et al.* Evidence of regio-specific glycosylation in human intestinal mucins:
912   presence of an acidic gradient along the intestinal tract. *J. Biol. Chem.* **278,** 46337-46348.
913   (2003).

914   22. Holmén Larsson, J. M., Thomsson, K. A., Rodríguez-Piñeiro, A. M., Karlsson, H. &
915   Hansson, G. C. Studies of mucus in mouse stomach, small intestine, and colon. III.
916   Gastrointestinal Muc5ac and Muc2 mucin O-glycan patterns reveal a regiospecific
917   distribution. *Am. J. Physiol. Gastrointest. Liver Physiol.* **305,** G357-363 (2013).

918   23. Larsson, J. M., Karlsson, H., Sjövall, H. & Hansson, G. C. A complex, but uniform O-
919   glycosylation of the human MUC2 mucin from colonic biopsies analyzed by nanoLC/MSn.
920   *Glycobiology* **19,** 756-766 (2009).

921   24. Moustafa, I. *et al.*  Sialic acid recognition by *Vibrio cholerae* neuraminidase. *J. Biol.*
922   *Chem.* **279,** 40819-40826 (2004).

923   25. Boraston, A. B., Ficko-Blean, E. & Healey, M. Carbohydrate recognition by a large
924   sialidase toxin from *Clostridium perfringens*. *Biochemistry (Mosc.)* **46,** 11352-11360 (2007).

925   26. Boraston, A. B., Bolam, D. N., Gilbert, H. J. & Davies, G. J. Carbohydrate-binding
926   modules: fine-tuning polysaccharide recognition. *Biochem. J.* **382,** 769-781 (2004).

927   27. Singh, A. K. *et al.* Unravelling the multiple functions of the architecturally intricate
928   *Streptococcus pneumoniae* β-galactosidase, BgaA. *PLoS Pathog.* **10,** e1004364 (2014).

929   28. Qin, J. *et al.* A human gut microbial gene catalog established by metagenomic
930   sequencing. *Nature* **464,** 59-65 (2010).

931   29. Crost, E. H., Tailford, L. E., Le Gall, G., Fons, M., Henrissat, B. & Juge, N. Utilisation of
932   mucin glycans by the human gut symbiont *Ruminococcus gnavus* is strain-dependent. *PloS*
933   *One* **8,** e76341 (2013).

934   30. Tailford, L. E. *et al.*  Discovery of intramolecular *trans*-sialidases in human gut microbiota
935   suggests novel mechanisms of mucosal adaptation. *Nat. Commun.* **6,** 7624 (2015).

936   31. Crost, E. H. *et al.* The mucin-degradation strategy of *Ruminococcus gnavus*: The
937   importance of intramolecular *trans*-sialidases. *Gut Microbes* **25,**1-11 (2016).

938   32. Ribeiro, J. P. *et al.* Characterization of a high-affinity sialic acid-specific CBM40 from
939   *Clostridium perfringens* and engineering of a divalent form. *Biochem J.* **473,** 2109-2118
940   (2016).

941   33. Luo, Y., Li, S. C., Chou, M. Y., Li, Y. T. & Luo, M., 1998. The crystal structure of an
942   intramolecular trans-sialidase with a NeuAc alpha2-->3Gal specificity. *Struct. Lond. Engl.* **6,**
943   521-530 (1993).

944   34. Connaris, H., Crocker, P. R. &Taylor, G. L. Enhancing the receptor affinity of the sialic
945   acid-binding domain of *Vibrio cholerae* sialidase through multivalency. *J. Biol. Chem.* **284,**
946   7339-7351 (2009).

947   35. Yang, L., Connaris, H., Potter, J. A., Taylor, G. L. Structural characterization of the
948   carbohydrate-binding module of NanA sialidase, a pneumococcal virulence factor. *BMC*
949   *Struct. Biol.* **15,** 15 (2015).

950 36. Xu, G., Potter, J. A., Russell, R. J., Oggioni, M. R., Andrew, P. W. & Taylor, G. L. Crystal
951 structure of the NanB sialidase from *Streptococcus pneumoniae*. *J. Mol. Biol.* **1384,** 436-449
952 (2008).

953 37. Owen, C. D., Lukacik, P., Potter, J. A., Sleator, O., Taylor, G. L. & Walsh, M. A.
954 *Streptococcus pneumoniae* NanC: Structural insights into the specificity and mechanism of a
955 sialidase that produces a sialidase inhibitor. *J. Biol. Chem.* **290,** 27736-27748 (2015).

956 38. Padler-Karavani, V. *et al.* Cross-comparison of protein recognition of sialic acid diversity
957 on two novel sialoglycan microarrays. *J. Biol. Chem.* **287,** 22593-22608 (2012).

958 39. Deng, L., Chen, X. & Varki, A. Exploration of sialic acid diversity and biology using
959 sialoglycan microarrays. *Biopolymers* **99,** 650-665 (2013).

960 40. Mayer, M. & Meyer, B. Characterization of ligand binding by saturation transfer
961 difference NMR spectroscopy. *Ang. Chem. Int. Ed.* **38,** 1784-1788 (1999).

962 41. Angulo, J. & Nieto, P. M. STD NMR: application to transient interactions between
963 biomolecules-a quantitative approach. *Eur. Biophys. J.* **40,** 1357-1369 (2011).

964 42. Mayer, M. & Meyer, B. Group epitope mapping by saturation transfer difference NMR to
965 identify segments of a ligand in direct contact with a protein receptor. *J. Am. Chem. Soc.*
966 **123,** 6108-6117 (2001).

967 43. Marchetti, R., *et al.* Rules of engagement" of protein–glycoconjugate interactions: a
968 molecular view achievable by using NMR spectroscopy and molecular modeling. *Chemistry*
969 *Open* **5,** 274-296 (2016).

970 44.Thomsson, K. A., Holmén-Larsson, J. M., Angström, J., Johansson, M. E., Xia L. &
971 Hansson, G. C. Detailed O-glycomics of the Muc2 mucin from colon of wild-type, core 1- and
972 core 3-transferase-deficient mice highlights differences compared with human MUC2.
973 *Glycobiology* **22,** 1128-39 (2012).

974 45. Homquist, L. & Ostman, B. The anomeric configuration of N-acetylneuraminic acid
975 released by the action of *Vibrio cholerae* neuraminidase. *FEBS Lett.* **60,** 327-330 (1975).

976 46. Pell G., Williamson M. P., Walters C., Du H., Gilbert H. J. & Bolam D. N. Importance of
977 hydrophobic and polar residues in ligand binding in the family 15 carbohydrate-binding
978 module from *Cellvibrio japonicus* Xyn10C. *Biochemistry* **42,** 9316-9323 (2003).

979 47. Abbott, D. W. & van Bueren, A. L. Using structure to inform carbohydrate binding module
980 function. *Curr. Opin. Struct. Biol.* **28,** 32-40 (2014).

981 48. Etzold, S. *et al.* Structural basis for adaptation of lactobacilli to gastrointestinal mucus.
982 *Environ. Microbiol.* **16,** 888-903 (2014).

983 49. Gunning, A. P., Kavanaugh, D., Thursby, E., Etzold, S., MacKenzie, D. A. & Juge, N.
984 Use of atomic force microscopy to study the multi-modular interaction of bacterial adhesins
985 to mucins. *Int. J. Mol. Sci.* **17,** pii: E1854 (2016).

986 50. Varki, N. M., Strobert, E., Dick, E. J. J., Benirschke, K. & Varki, A. Biomedical differences
987 between human and nonhuman hominids: potential roles for uniquely human aspects of
988 sialic acid biology. *Annu. Rev. Pathol.* **6,** 365-393 (2011).

989 51. Tangvoranuntakul, P., *et al.* Human uptake and incorporation of an immunogenic
990 nonhuman dietary sialic acid. *Proc. Natl. Acad. Sci. U.S.A.* **100,** 12045-12050 (2003).

991 52. Ficko-Blean, E. & Boraston, A. B. Insights into the recognition of the human glycome by
992 microbial carbohydrate-binding modules. *Curr. Opin. Struct. Biol.* **22,** 570-577 (2012).

993 53. Hervé, C., Rogowski, A., Blake, A. W., Marcus, S. E., Gilbert, H. J. & Knox, J. P.
994 Carbohydrate-binding modules promote the enzymatic deconstruction of intact plant cell
995 walls by targeting and proximity effects. *Proc. Natl. Acad. Sci. U. S. A.* **107,** 15293–15298.
996 (2010).

997  54. El Kaoutari, A., Armougom, F. Gordon, J. I., Raoult, D. & Henrissat, B. The abundance
998  and variety of carbohydrate-active enzymes in the human gut microbiota. *Nat. Rev.*
999  *Microbiol.* **11,** 497-504 (2013).

1000  55. Yu, H. *et al.* Sequential one-pot multienzyme chemoenzymatic synthesis of
1001  glycosphingolipid glycans. *J. Org. Chem.* **81,** 10809-10824 (2016).

1002  56. Liu, H. & Naismith, J. H. A simple and efficient expression and purification system using
1003  two newly constructed vectors. *Protein Expr. Purif.* **63,** 102-111 (2009).

1004  57. Bergfors, T. Seeds to crystals. *J. Struct. Biol.* **142,** 66-76 (2003).

1005  58. Otwinowski, Z. & Minor, W. Processing of X-ray diffraction data collected in oscillation
1006  mode, in: Macromolecular crystallography, Part A, Methods in enzymology. Academic Press,
1007  New York, pp. 307–326 (1997).

1008  59. McCoy, A.J., Grosse-Kunstleve, R.W., Adams, P.D., Winn, M.D., Storoni, L.C. & Read,
1009  R.J. Phaser crystallographic software. *J. Appl. Crystallogr.* **40,** 658-674 (2007).

1010  60. Winn, M. D. *et al.* Overview of the CCP 4 suite and current developments. *Acta*
1011  *Crystallogr. D Biol. Crystallogr.* **67,** 235–242 (2011).

1012  61. Murshudov, G. N. *et al.* REFMAC5 for the refinement of macromolecular crystal
1013  structures. *Acta Crystallogr. D Biol. Crystallogr.* **67,** 355-367 (2011).

1014  62. Emsley, P., Lohkamp, B., Scott, W. G. & Cowtan, K. Features and development of Coot.
1015  *Acta Crystallogr. D Biol. Crystallogr.* **66,** 486-501 (2010).

1016  63. Joosten, R. P., Joosten, K., Murshudov, G. N. & Perrakis, A. PDB_REDO: constructive
1017  validation, more than just looking for errors. *Acta Crystallogr. D Biol. Crystallogr.* **68,** 484-496
1018  (2012).

1019  64. Chen, V. B. *et al.* MolProbity: all-atom structure validation for macromolecular
1020  crystallography. *Acta Crystallogr. D Biol. Crystallogr.* **66,** 12-21 (2010).

1021  65. Karplus, P. A. & Diederichs, K. Linking Crystallographic Model and Data Quality. *Science*
1022  **336,** 1030-1033 (2012).

1023  66. Finn, R. D. *et al.* The Pfam protein families database: towards a more sustainable future.
1024  *Nucleic Acids Res.* **44,** D279-285 (2016).

1025  67. Felsenstein, J. PHYLIP - Phylogeny Inference Package (Version 3.2). *Cladistics* **5,** 164-
1026  166 (1989).

1027  68. Rice, P. Longden, I. & Bleasby, A. EMBOSS: The European Molecular Biology Open
1028  Software Suite. *Trends Genet.* **16,** 276-277 (2000).

1029  69. Liquori, G. E. *et al.* In situ characterization of O-linked glycans of Muc2 in mouse colon.
1030  *Acta Histochem.* **114,** 723-732 (2012).

1031

29

1039 Carmen Pin for her help with statistical analyses and Emmanuelle Crost for *R. gnavus*

1040 anaerobic culture.

1041

1042 **Author Contributions:** NJ conceived the study and wrote the manuscript with contribution

1043 from all co-authors. CDO carried out sub-cloning, produced the proteins (*Rg*CBM40,

1044 *Rg*GH33, *Rg*NanH) and solved CBM40 crystal structures under GLT's supervision. LET

1045 carried out the cloning, heterologous expression, mutagenesis and CD analysis of proteins

1046 (*Rg*CBM40, *Rg*GH33, *Rg*NanH) and carried out binding assays (ITC and ELISA) and

1047 enzyme kinetics (HPAEC), TS and LV carried out the immuno- histo/cytochemistry

1048 experiments, ST purified the mucins from human cell lines and mouse models, KL

1049 characterized the glycosylation profile of mucins by MS, MH contributed to the production of

1050 *Rg*CBM40, *Rg*NanH and *Rg*GH33, and to CD and ELISA experiments. RL contributed to the

1051 production of *Rg*CBM40, and to the CD and ITC experiments, AB performed the western

1052 blot analysis and prepared cells for TEM.  AB, KL, LET, LV, ST, and TS worked under NJ's

1053 supervision, MH and RL worked under LET's supervision. JW performed the bioinformatics

1054 analyses. SM carried out the STD NMR experiments under JA's supervision, ZK performed

1055 the glycan microarray screening under AV's supervision, HY synthesized some of the

1056 sialosides used in this study under XC's supervision.

1057
1058 **Conflict of interest:** The authors declare no conflict of interest.
1059
1060
1061

**Figure Legends**
1063

1064 **Figure 1 Crystal structure of *Rg*CBM40 in complex with .3'SL and 6'SL** (**a**) *Rg*CBM40 is

1065 shown in a cartoon representation with a rotation of 90° around the *x* axis. (**b**) The protein

1066 crystallised as a dimer with the ligand binding site at the dimer interface. The binding sites

1067 are shown occupied by 6'SL trisaccharides (Neu5Ac: cyan, galactose: blue, glucose:

1068 orange). (**c**) SEC-MALS performed with full length *Rg*NanH (77 kDa). The SEC-MALS

1069 predicted molecular weight was 73 kDa, indicating that *Rg*NanH is monomeric in solution.

1070 Bound 3'SL (**d**) and 6'SL (**e**) are shown with their corresponding Fo-Fc omit maps at 2 σ

1071 (light cyan), 3 σ (orange), and 5 σ (magenta). The omit maps are carved at 1.6 Å around the

1072 bound ligand. For 3'SL, the map is carved around a dummy glucose residue to indicate the

1073 presence of partial electron density. A close-up view of *Rg*CBM40 binding site is shown with

1074 (**f**) 3'SL and (**g**) 6'SL. The Neu5Ac residue is shown in cyan and the galactose residue as

1075 black lines, for clarity the glucose residue is not shown. Interacting *Rg*CBM40 residues are

1076 shown in green with black dashed lines indicating hydrogen bonding interactions. A semi-

1077 transparent surface indicates the hydrophobic surface.

1078

1079 **Figure 2. CBM40 structural alignment**

1080 Structure-based alignment (α-helices and β-strands respectively in red and yellow) of

1081 CBM40 domains of *Rg*CBM40 with *C. perfringens Cp*CBM40_NanJ (PDB code 2V73) and

1082 *Cp*CBM40_NanI (PDB code 5FRA), *M. decora Md*CBM40_NanL (PDB code 1SLI) and *S.*

1083 *pneumoniae Sp*CBM40_NanA (PDB code 4C1W), *Sp*CBM40_NanB (PDB code 2VW0) and

1084 *Sp*CBM40_NanC (PDB code 4YZ5) and *Vc*CBM40_NanH structure (PDB code 2W68).

1085 Amino acids identified as binding sites are highlighted in blue. *Rg*CBM40 residues Ile95,

1086 Asp110, Tyr116, Glu126, Arg128, Arg204 and Tyr210 are at positions 104, 119, 125, 135,

1087 137, 226 and 233 of the alignment. The alignment supplemented with other canonical and

1088 *Vibrio*-type CBM40 sequences, used to create the pHMM using HMMER3, is shown in

1089 **Supplementary Fig. 4**.

1090

1091 **Figure 3. Distance-based tree of canonical and *Vibrio*-type CBM40 sequences**

1092 Tree of 51 non-redundant sequences (80% identity level) calculated by neighbour-joining

1093 using evolutionary distances estimated by applying the PMB model of amino acid changes,

1094 including all sites and using a uniform rate of evolution. The representative sequences

1095 corresponding most closely (at least 97% identical) to the 7 bacterial structure-determined

1096 sequences are shown with symbols, coloured in accordance with **Supplementary Fig. 1**:

1097 "A", *Sp*CBM40_NanA; "B", *Sp*CBM40_NanB "C", *Sp*CBM40_NanC; "I", *Cp*CBM40_NanI;

1098 "J", *Cp*CBM40_NanJ; "*R*", *Rg*CBM40; "*V*", *Vc*CBM40_NanH. Additionally, "L" denotes

1099     *Md*CBM40_NanL closest to the bacterial sequence of highest identity (70% identical to

1100     *Rg*CBM40) as only bacterial sequences were searched.

1101

1102     **Figure 4. Sialoglycan microarray analysis of binding specificities of *Rg*CBM40 and**

1103     ***Rg*GH33 D282A.** Binding of the recombinant proteins *Rg*CBM40 and *Rg*GH33 D282A at 20

1104     and 200 µg mL$^{-1}$, respectively are presented (n=4, SD). Heat map was generated using the

1105     method as previously described[38,39]. Binding was ranked as (glycan average RFU/ maximum

1106     glycan average RFU)*100. Red and white represent the maximum and minimum,

1107     respectively. R1 represents propyl-azide as the spacer.

1108

1109     **Figure 5. STD NMR analysis of *Rg*CBM40 binding to sialoglycans,** (**a**) Reference (top)

1110     and difference (bottom) spectra of 3'SL and 6'SL. The strongest signals from the Neu5Ac's

1111     protons are labelled in the difference spectra. (**b**) Binding epitope mapping from STD NMR

1112     of 3'SL and 6'SL. Legend indicates relative STD intensities normalised at H7: blue, 0–24%;

1113     yellow, 25–50%, red 51–100%; larger red dots indicate values over 100%. Sialic acid is the

1114     main recognition element. (**c**) Binding epitopes mapping from STD NMR of Neu5Gcα2-3Lac

1115     and Neu5Gcα2-6Lac. Legend as above. Sialic acid is the main recognition element. The

1116     strongest STD intensities from CH2 and the H3s, suggest a reorientation of the Neu5Gc ring

1117     in the binding pocket, in comparison to 3'SL and 6'SL.

1118

1119     **Figure 6.  ITC isotherms of *Rg*CBM40 to sialoglycans. (a)** *Rg*CBM40 binding to 3'SL, **(b)**

1120     *Rg*CBM40 binding to 6'SL, **(c)** *Rg*CBM40 binding to 3'SLGc, **(d)** *Rg*CBM40 binding to

1121     Neu5Ac. The Kd is indicated in mM. *This value is an estimate as the Kd is too high to

1122     determine with the concentration of sugar used.

1123

1124     **Figure 7. ELISA of *Rg*CBM40 binding to purified mucins.**

1125     (**a**) *Rg*CBM40 binding to a range of purified mucins; mucin 2 (MUC2) and mixed mucins

1126     (mucins) from human cell line LS174T, purified pig gastric mucin (pPGM), and murine

1127     mucins from germ free (GF), wild type (WT), and *C3GnT*$^{-/-}$ mice. (**b**) Correlation of

1128     *Rg*CBM40 binding with % sialylated structure for each mucin tested. *Rg*CBM40 was

1129     incubated with immobilised mucins and binding determined by ELISA. The % sialylated

1130     structures was determined by MS. (**c**) *Rg*CBM40 binding to LS174T MUC2 which has been

1131     treated chemically (TFA) or enzymatically with a sialidase from *Clostridium perfringens* (*Cp*),

1132     *Salmonella typhimurium* (*St*), *Akkermansia muciniphila* (*Am*) or *Ruminococcus gnavus* (*Rg*)

1133     **(d)** *Rg*CBM40 binding to LS174T MUC2 in competition with sugars. *Rg*CBM40 has been

1134     preincubated with the indicated sugars. In all cases, *Rg*CBM40 was incubated with

1135     immobilised mucins and visualised by an ELISA using anti-sialidase primary antibody and an

1136 anti-rabbit secondary antibody conjugated to horseradish peroxidase. The enzyme was

1137 incubated with TMB and the absorbance at 450 nm (A450) measured. The error bars show

1138 the standard error of the mean (SEM) of three replicates. P values are indicated; NS-not

1139 significant, *-p<0.05, **-p<0.005, ***-p<0.005.

1140

1141 **Figure 8. *Rg*CBM40 binding to mucus-producing cells and intestinal tissue sections.**

1142 (**a**) Immunostaining pattern for *Rg*CBM40 on LS174T cells correlated with mucin (MUC2)

1143 and lectin (SNA) staining, all shown in green. No staining was observed in *Rg*CBM40-free

1144 sample (Blank). (**b**) Immunostaining pattern for *Rg*CBM40 on cryosections of mouse colon

1145 correlated with mucin (Muc2) and lectin (SNA) staining, all shown in green. No staining was

1146 observed in *Rg*CBM40-free sample (Blank). Cell nuclei were counterstained with DAPI,

1147 shown in blue. (**c**) Sialidase pre-treatment of mouse colonic cryosections markedly reduced

1148 the binding of *Rg*CBM40 and SNA lectin. Cell nuclei were counterstained with DAPI, shown

1149 in blue. (**d**) *Rg*CBM40 competition assay with SNA on cryosections of mouse colon.

1150 *Rg*CBM40 is shown in green. Cell nuclei were counterstained with DAPI, shown in blue. No

1151 *Rg*CBM40 specific staining was detectable when SNA was present. (**e**) *R. gnavus* binding

1152 competition assay with SNA on cryosections of mouse colon. *R. gnavus* ATCC 29149 was

1153 incubated on sequential cryosections of mouse colon with or without SNA treatment and is

1154 shown in red. The mucus layer is shown in green. Sequential sections were required as both

1155 antibodies were raised in the same species. Cell nuclei were counterstained with DAPI,

1156 shown in blue.  No *R. gnavus* staining was detectable when SNA was present. Appropriate

1157 primary antibody and secondary antibody only controls are also shown underneath each

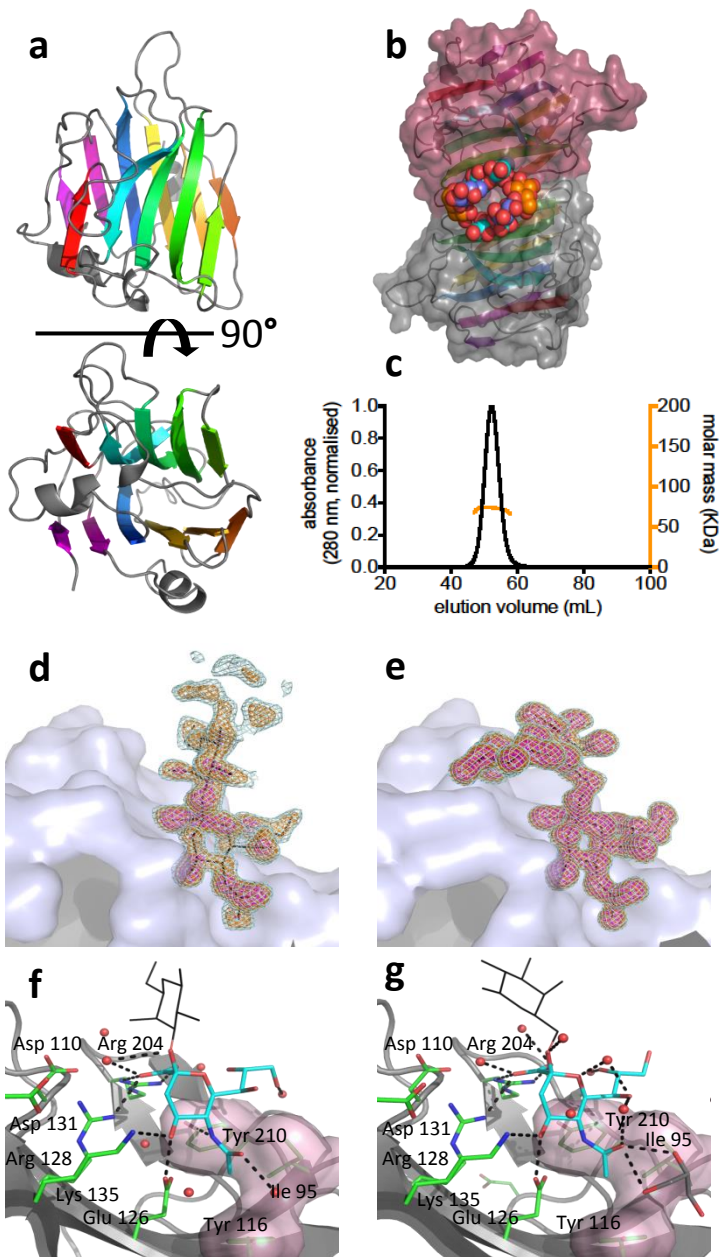1158 panel, showing some background staining. Scale bar: 20 µm.

1159

1160

1161

1162 **Table 1.** Data collection and refinement statistics. Values in parentheses refer to the highest

1163 resolution shell. For the 3'SL and 6'SL complexes the data was over 90% complete to a

1164 resolution of 1.85 Å and 1.56 Å respectively.

1165

| Dataset | Apo | 3'SL | 6'SL |
|---|---|---|---|
| Data collection | | | |
| Spacegroup | P21 | P21 | P21 |
| Cell dimensions | | | |
| a, b, c (Å) | 46.7, 72.8, 51.3 | 48.8, 72.4 51.5 | 48.7, 72.2, 51.4, |
| $\beta$ (°) | 104.9 | 105.1 | 103.9 |
| Resolution | 50 – 1.73 (1.76 – 1.73) | 39.48 – 1.37 (1.41 – 1.37) | 49.91 – 1.30 (1.34 – 1.30) |
| $R_{merge}$ | 0.03 (0.14) | 0.04 (0.34) | 0.03 (0.15) |
| $I / \sigma I$ | 47.3 (9.6) | 22.9 (3.0) | 32.0 (4.9) |
| Completeness | 91.8 (51.3) | 74.5 (11.2) | 83.9 (13.6) |
| Redundancy | 3.7 (2.4) | 4.3 (2.4) | 5.1 (1.4) |
| | | | |
| Refinement | | | |
| Resolution | 50 – 1.73 (1.76 – 1.73) | 39.48 – 1.37 (1.41 – 1.37) | 49.91 – 1.30 (1.34 – 1.30) |
| No. reflections | 31570 | 51221 | 67097 |
| $R_{work}$ / $R_{free}$ | 0.160/0.194 (0.82) | 0.152/0.187 (0.81) | 0.134/0.154 (0.87) |
| No. of atoms | 3145 | 3424 | 3704 |
| Protein | 2807 | 2850 | 3076 |
| Ligand | 0 | 81 | 123 |
| Water | 338 | 508 | 527 |
| *B*-factors | | | |
| Protein | 19.4 | 16.6 | 10.4 |
| Ligand/ion | | 36.4 | 21.7 |
| Water | 28.1 | 31.7 | 27.3 |
| R.m.s.d | 0.011 | 0.012 | 0.015 |
| Bond lengths (Å) | 0.011 | 0.012 | 0.015 |
| Bond angle (°) | 1.55 | 1.66 | 1.77 |

1166 *Values in parentheses are for the highest-resolution shell. One crystal was used for each
1167 structure.
1168
1169
1172

a

b

c

d

e

f

Asp 110    Arg 204

Asp 131

Arg 128                          Tyr 210

Lys 135                           Ile 95

Glu 126         Tyr 116

g

Asp 110    Arg 204

Asp 131

Arg 128                          Tyr 210
                                 Ile 95
Lys 135
Glu 126         Tyr 116

Multiple sequence alignment

```
              10        20        30        40        50        60        70
SpCBM40_NanA  ---------VE--TVIEKEDVET-NASNGQ----------------------------------RVDL--SS-EL
SpCBM40_NanB  --------------IFQGGSYQLNN-K-S-----------------------------------IDI--SSLLL
SpCBM40_NanC  ----------PVLEKNNVTLTG-G-G-------------------------------------ENV--TKELK
CpCBM40_NanI  ---SPDPNWELLSSLGEYKDINL-ESSNA---------------------------------SNI--TY-DL
CpCBM40_NanJ  LNVYEIKGEVD--EIANYGNLKITKEEER---------------------------------VNI--TG-DL
McCBM40_NanL  -------PEG--ILMEKNNVDI-AEGQG----------------------------------YSLDQEA-GA
RgCBM40       ---------SV--PVLQKEGIEI-SEGTG---------------------------------YDLSKEP-GA
VcCBM40_NanH  ---------S-NAALFDYN---ATGDTEFDSPAKQGWMQDNTNNGSGVLTNADGMPAWLVQGIGGRAQWTYSL--STNQH

              90        100       110       120       130       140       150
SpCBM40_NanA  DKLKKLENATVHMEFKPDAKAPAFYNLFSVSSATK--KDEYFTMAVYNN-TATLESRGS--DGKQ---FYN-NYNDAPLK
SpCBM40_NanB  DKL-SGESQTVVMKFKADKP-NSLQALFGLSNSKAGFKNNYFSIFMRDSGEIGVEIRD----AQK---GIN-YLFSRPAS
SpCBM40_NanC  DKFTSGD-FTVVIKYNQSSE-KGLQALFGISNSKPGQQNSYVDVFLRDNGELGMEARD--TSSN-----KN-NLVSRPAS
CpCBM40_NanI  EKYKNLDEGTIVVRFNS-KD-SKIQSLLGISNSK--TKNGYFNFYVTNS-RVGFEIRNQKNEGNTQNGTENLVHMYKDVA
CpCBM40_NanJ  EKFSSLEEGTIVTRFNM-ND-TSIQSLIGLSDGNK--ANNYFSLYVS-GGKVGYEIRR--QEGNG---DFN-VHHSADVT
McCBM40_NanL  KYVKAMTQGTIILSYKSTSE-NGIQSLFSVGNSTAGNQDRHFHIYITNSGGIGIELRN--TDGV-----FN-YTLDFPAS
RgCBM40       ATVKALEQGTIVISYKTTSE-NAIQSLLSVGNGTKGNQDRHFHLYITNAGGVGMELRN--TDGE-----FK-YTLDCPAA
VcCBM40_NanH  AQASS-FGWRMTTEMKVLSG---GMITNYYANG---TQRVLPIISLDSSGNLVVEF----EGQT---GR---TVLATGT

              170       180       190       200       210       220       230
SpCBM40_NanA  VKP---GQW--NSVTFTVEKPTAELPKGRVRLYVNGVLSRT-SLRSGNFIKDMPDVTHVQIGATKRAN-NTYW-GSNLQI
SpCBM40_NanB  LWGKHKGQAVENTLVFVSDSKDK-----TYTMYVNGIEVFSETVDTFLPISNINGIDKATLGAVNRE-GKEHY-LAKGSI
SpCBM40_NanC  VWGKYKQEAVTNTVAVVADSVKK-----TYSLYANGTKVVEKKVDNFLNIKDIKGIDYYMLGGVKRA-GKTAF-GFNGTL
CpCBM40_NanI  L----NDGD--NTVALKIEKN------KGYKLFLNGKMIKEVKDTNTKFLNNIENLDSAFIGKTNRYGQSNEY-NFKGNI
CpCBM40_NanJ  F----NRGI--NTLALKIEK------GIGAKIFLNGSLVKTVSDPNIKFLNAI-NLNSGFIGKTDRANGYNEY-LFRGNI
McCBM40_NanL  VRALYKGERVFNTVALKADAANK-----QCRLFANGELLATLDKDAFKFISDITGVDNVTLGGTKRQ-GKIAY-PFGGTI
RgCBM40       VRGSYKGERVSNTVALKADKENK-----QYKLFANGELIATLDQEAFKFISDITGVDNVMLGGTMRQ-GTVAY-PFGGSI
VcCBM40_NanH  AATEY------HKFELVFLPGSNP----SASFYPDGKLIRDNIQPTA------SKQNMIVWGN------GSSNTDGVAAY

              250       260
SpCBM40_NanA  RNLTVYNRALTPEEVQKRSQLFK
SpCBM40_NanB  DEISLFNKAISDQEVSTIPLSNP
SpCBM40_NanC  ENIKFFNSALDEETVKKMTTNA-
CpCBM40_NanI  GFMNIYNEPLGDDYLLSKTGETK
CpCBM40_NanJ  DFMNIYDKPVSDNYLLRKTGETK
McCBM40_NanL  GDIKVYSNALSDEELIQATGVTT
RgCBM40       ERMQVYRDVLSDDELIAVTGKT-
VcCBM40_NanH  RDIKFEIQGD-------------
```

| Glycan Structure | *Rg*CBM40 | *Rg*GH33 D282A | Rank |
|---|---|---|---|
| Neu5Acα6GalNAcαR1 | | | 100 |
| Neu5Acα3Galβ4GlcNAcβR1 | | | 50 |
| Neu5Acα3Galβ3GlcNAcβR1 | | | 0 |
| Neu5Acα3Galβ3GalNAcαR1 | | | |
| Neu5Acα6Galβ4GlcNAcβR1 | | | |
| Neu5Acα6Galβ4GlcβR1 | | | |
| Neu5Acα3Galβ4GlcβR1 | | | |
| Neu5Acα3GalβR1 | | | |
| Neu5Acα6GalβR1 | | | |
| Neu5Acα3Galβ3GalNAcβR1 | | | |
| Neu5Acα8Neu5Acα3Galβ4GlcβR1 | | | |
| Neu5Acα8Neu5Acα8Neu5Acα3Galβ4GlcβR1 | | | |
| Neu5Acα3Galβ4(Fucα3)GlcNAcβR1 | | | |
| Neu5Acα3Galβ4(Fucα3)GlcNAc6SβR1 | | | |
| Neu5Acα3Galβ3GlcNAcβ3Galβ4GlcβR1 | | | |
| Neu5Acα3Galβ4GlcNAc6SβR1 | | | |
| Neu5Acα6(Neu5Acα3)Galβ4GlcβR1 | | | |
| Neu5Acα6(Neu5Gcα3)Galβ4GlcβR1 | | | |
| Neu5Acα6(Kdnα3)Galβ4GlcβR1 | | | |
| Neu5Acα8Neu5Gcα3Galβ4GlcβR1 | | | |
| Neu5Acα8Neu5Gcα6Galβ4GlcβR1 | | | |
| Neu5Acα8Neu5Acα6Galβ4GlcβR1 | | | |
| Neu5,9Ac₂α3Galβ4GlcNAcβR1 | | | |
| Neu5,9Ac₂α6Galβ4GlcNAcβR1 | | | |
| Neu5,9Ac₂α3Galβ3GlcNAcβR1 | | | |
| Neu5,9Ac₂α3Galβ3GalNAcαR1 | | | |
| Neu5,9Ac₂α6GalNAcαR1 | | | |
| Neu5,9Ac₂α3GalβR1 | | | |
| Neu5,9Ac₂α6GalβR1 | | | |
| Neu5,9Ac₂α3Galβ3GalNAcβR1 | | | |
| Neu5,9Ac₂α6Galβ4GlcβR1 | | | |
| Neu5,9Ac₂α3Galβ4GlcβR1 | | | |
| Neu5Gcα6GalNAcαR1 | | | |
| Neu5Gcα3Galβ4GlcNAcβR1 | | | |
| Neu5Gcα3Galβ3GlcNAcβR1 | | | |
| Neu5Gcα3Galβ3GalNAcαR1 | | | |
| Neu5Gcα6Galβ4GlcNAcβR1 | | | |
| Neu5Gcα6Galβ4GlcβR1 | | | |
| Neu5Gcα3Galβ4GlcβR1 | | | |
| Neu5Gcα3GalβR1 | | | |
| Neu5Gcα6GalβR1 | | | |
| Neu5Gcα3Galβ3GalNAcβR1 | | | |
| Neu5Gcα3Galβ4(Fucα3)GlcNAcβR1 | | | |
| Neu5Gcα3Galβ4(Fucα3)GlcNAc6SβR1 | | | |
| Neu5Gcα3Galβ3GlcNAcβ3Galβ4GlcβR1 | | | |
| Neu5Gcα3Galβ4GlcNAc6SβR1 | | | |
| Neu5Gcα8Neu5Acα3Galβ4GlcβR1 | | | |
| Neu5Gcα8Neu5Gcα3Galβ4GlcβR1 | | | |
| Neu5Gc9Acα3Galβ4GlcNAcβR1 | | | |
| Neu5Gc9Acα6Galβ4GlcNAcβR1 | | | |
| Neu5Gc9Acα3Galβ3GlcNAcβR1 | | | |
| Neu5Gc9Acα3Galβ3GalNAcαR1 | | | |
| Neu5Gc9Acα6GalNAcαR1 | | | |
| Neu5Gc9Acα3GalβR1 | | | |
| Neu5Gc9Acα6GalβR1 | | | |
| Neu5Gc9Acα3Galβ3GalNAcβR1 | | | |
| Neu5Gc9Ac6Galβ4GlcβR1 | | | |
| Neu5Gc9Ac3Galβ4GlcβR1 | | | |
| Kdnα8Neu5Acα3Galβ4GlcβR1 | | | |
| Kdnα8Neu5Gcα3Galβ4GlcβR1 | | | |

**(a)** Neu5Acα2-3Lac / Neu5Acα2-6Lac

H8 H9' H5 H6 H4 H7 Me H3e H3a

**(b)** Strong ● Medium ● Weak ●

**(c)** Strong ● Medium ● Weak ●

**a** Blank | *Rg*CBM40 | SNA | MUC2

**b** Blank | *Rg*CBM40 | SNA | Muc2

**c** *Rg*CBM40 (no sialidase) | *Rg*CBM40 (sialidase) | SNA (no sialidase) | SNA (sialidase)

**d** *Rg*CBM40 only | *Rg*CBM40 merged | *Rg*CBM40 and SNA | *Rg*CBM40 and SNA merged

**e** *R. gnavus* ATCC29149 only | Muc2 | *R. gnavus* and SNA | Muc2

Primary antibody control | Primary antibody control | Secondary antibody control | Secondary antibody control