

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26

Targeted protein evolution in the gut microbiome by diversity-generating retroelements

Benjamin R. Macadangdang^{1,2,*}, Yanling Wang³, Cora Woodward², Jessica I. Revilla⁴, Bennett M. Shaw⁵, Kayvan Sasaninia³, Sara K. Makanani^{3,6}, Chiara Berruto³, Umesh Ahuja^{2,*}, Jeff F. Miller^{2,3,6,8,*}

¹Division of Neonatology and Developmental Biology, Department of Pediatrics, David Geffen School of Medicine at the University of California, Los Angeles, Los Angeles, CA, United States

²California NanoSystems Institute, Los Angeles, CA, United States

³Department of Microbiology, Immunology and Molecular Genetics, University of California, Los Angeles, California, United States

⁴Department of Molecular, Cell, and Developmental Biology, University of California, Los Angeles, Los Angeles, CA, United States

⁵David Geffen School of Medicine, University of California, Los Angeles, Los Angeles, CA, United States

⁶Molecular Biology Institute, University of California, Los Angeles, Los Angeles, CA, United States

⁷Division of Biology and Biological Engineering, California Institute of Technology, Pasadena, CA, United States

⁸Lead contact

*Correspondence: bmacadangdang@mednet.ucla.edu (BRM), uahuja@ucla.edu (UA), jfmiller@g.ucla.edu (JFM)

27 **Summary**

28 Diversity-generating retroelements (DGRs) accelerate evolution by rapidly diversifying variable proteins.
29 The human gastrointestinal microbiota harbors the greatest density of DGRs known in nature, suggesting
30 they play adaptive roles in this environment. We identified >1,100 unique DGRs among human-
31 associated *Bacteroides* species and discovered a subset that diversify adhesive components of Type V
32 pili and related proteins. We show that *Bacteroides* DGRs are horizontally transferred across species,
33 that some are highly active while others are tightly controlled, and that they preferentially alter the
34 functional characteristics of ligand-binding residues on adhesive organelles. Specific variable protein
35 sequences are enriched when *Bacteroides* strains compete with other commensal bacteria in gnotobiotic
36 mice. Analysis of >2,700 DGRs from diverse phyla in mother-infant pairs shows that *Bacteroides* DGRs
37 are preferentially transferred to vaginally delivered infants where they actively diversify. Our observations
38 provide a foundation for understanding the roles of stochastic, targeted genome plasticity in shaping host-
39 associated microbial communities.

40 Introduction

41 Natural selection acts on preexisting genetic variation to favor adaptive phenotypes. In bacteria, this
42 variation primarily arises from mutations and horizontal gene transfer¹. Nonetheless, genomic integrity is
43 vital for survival and reproductive success, prompting the evolution of mechanisms to channel
44 mutagenesis to localized hotspots within genomes²⁻⁶ and to limit horizontal transfer⁷⁻⁹. Of the known
45 systems that target mutagenesis, diversity-generating retroelements (DGRs) found in bacteria, archaea,
46 and their viruses can generate some of the most extensive repertoires of DNA and protein sequence
47 variants observed in nature^{10,11}. Through hypermutation of genes via a mechanism termed mutagenic
48 retrohoming, a single DGR can produce up to 10^{30} unique variable protein sequences¹². DGR target
49 genes often encode ligand-binding proteins with mutations strategically confined to a discrete subset of
50 codons within a variable repeat sequence (VR) that participates in ligand interactions¹³. The remaining
51 codons in VR, many of which encode structural scaffold residues, remain unmodified^{14,15}. This pattern of
52 mutagenesis rapidly diversifies protein function without disrupting structure, leading to readily evolvable
53 ligand-binding capabilities¹⁰.

54

55 In addition to a diversified target gene that encodes a variable protein, a typical DGR includes a template
56 repeat (TR) which is similar but not identical to VR, a uniquely promiscuous reverse transcriptase (RT),
57 and one or more accessory genes (Figure 1A). During mutagenic retrohoming, an RNA intermediate
58 encoded by TR functions as a substrate for reverse transcription by the DGR RT, which selectively
59 mismatches adenine residues. This results in random incorporation of any of the four nucleotides into a
60 cDNA molecule at positions corresponding to TR-adenines. Adenine-mutagenized cDNA is then
61 integrated into VR, replacing the parental allele¹⁶⁻¹⁸. Because TR is unaltered during this process, and
62 all *cis*- and *trans*-acting factors required for mutagenic retrohoming remain intact, repeated rounds of VR
63 diversification can occur indefinitely to optimize variable protein function.

64

65 The human gastrointestinal (GI) microbiome is enriched in DGRs to an extent that exceeds any other
66 ecosystem characterized to date¹⁹, yet the dynamics and roles of accelerated protein evolution in this
67 environment have never been systematically interrogated²⁰. Strains belonging to the Bacillota phylum
68 and the Fibrobacteres, Chlorobi, and Bacteroidota (FCB) superphylum harbor the vast majority of human
69 microbiome-associated DGRs. As constituents of the FCB group, *Bacteroides* spp. are prominent
70 members of the GI microbiome, forming long-term associations with their human host²¹. They impart
71 health benefits such as secretion of anti-inflammatory molecules^{22,23} and short chain fatty acids
72 (SCFAs)²⁴, but are also known to include pathobionts²⁵. Their amenability to *in vitro* cultivation and

73 genetic manipulation positions them as a clinically pertinent model for exploring DGR dynamics and
74 function in the context of their natural environment.

75

76 Single-nucleotide substitutions such as those produced by DGRs, as well as larger genomic variants
77 such as indels²⁶ or the presence of mobile genetic elements (MGEs)^{27,28}, can confer characteristics that
78 have profound effects on bacterial phenotypes and fitness. Conventional metagenomic methods
79 commonly lose variant information during assembly and may encounter challenges when categorizing
80 mobile elements during binning²⁹⁻⁴⁷. Thus, mechanisms underlying accelerated evolution, such as
81 horizontally transferable, DGR-driven mutagenic retrohoming, often evade detection by conventional
82 pipelines. Understanding the functions and phenotypes of complex microbial communities can best be
83 accomplished in the context of extant genotypic variation. To that end, we explored how rapid, targeted
84 evolution by DGRs contributes to *Bacteroides* diversity in the GI tract to gain a foundational understanding
85 of how genome plasticity molds host-associated microbial communities.

86

87 Here, we present a systematic analysis of DGR-mediated accelerated evolution in *Bacteroides* species.
88 By analyzing over 1,100 reference genomes, we found that DGRs are prevalent in *Bacteroides* and that
89 a large class of diversified variable proteins have homology to Type V pilins⁴⁸. We show that *Bacteroides*
90 DGRs can be transferred between strains, providing a mechanism for the horizontal transfer of
91 accelerated evolvability. Some *Bacteroides* DGRs are highly active *in vitro* and *in vivo*, and host-encoded
92 factors can further modulate mutagenic retrohoming. In the presence of competition with non-*Bacteroides*
93 strains *in vivo*, diversified *Bacteroides* VR regions can converge to encode similar protein sequences
94 despite being comprised of unique DNA sequences. Finally, by analyzing metagenomic datasets derived
95 from mother-infant pairs, we show that Bacteroidota DGRs are preferentially passed between mothers
96 and infants, where nearly 75% of transferred DGRs adopt a new, predominating VR haplotype. Our
97 results demonstrate that *Bacteroides* DGRs evolve bacterial proteins in the GI microbiome and are active
98 during periods of community instability.

99

100 **RESULTS**

101

102 **DGRs are widespread and diverse in *Bacteroides* species**

103 We analyzed a set of 1,103 *Bacteroides* reference genomes encompassing 47 species from the NCBI
104 RefSeq database⁴⁹ to determine the distribution of DGRs across taxa, to understand their evolutionary
105 relationships and modes of transmission, and to identify functional motifs in the diversified proteins
106 (Supplemental Figure 1A). From this dataset, we found 1,113 unique DGRs distributed across 618

107 *Bacteroides* isolates (Figure 1B-D). These represented 29 of the 47 species in our dataset and
108 encompassed 11 of the most abundant species in humans (Supplemental Table 4). A solitary DGR was
109 found in 340 isolates (31%), while an additional 278 isolates (25%) contained multiple DGRs, with some
110 genomes of *B. acidifaciens*, *B. xylanisolvens*, and *B. ovatus* harboring up to five unique elements (Figure
111 1B, Supplemental Figure 1B). DGRs were enriched in these and other species while *B. fragilis* and *B.*
112 *intestinalis* maintained fewer elements than average (Figure 1C). These results show that DGRs are
113 prevalent and abundant within *Bacteroides* spp. compared to most other bacterial taxa^{19,50}.

114

115 To gain insight into the relationships between *Bacteroides* DGRs, we built a phylogenetic tree using DGR
116 RT sequences (Figure 1D). Non-DGR RTs identified within the same genome set were included for
117 comparison⁵¹. We also inspected adjacent sequences for loci related to prophages⁵², integrative and
118 conjugative elements (ICEs)⁵³, or plasmids to identify DGRs that are likely to reside on MGEs^{19,20}. In the
119 absence of such evidence, DGRs were classified as "cellular" to indicate their presence in cellular
120 genomes. Insights into the potential adaptive roles were obtained by clustering variable proteins and
121 searching for similarities to protein domains of known function⁵⁴. Clusters were further categorized into
122 larger groups based on shared domains with the greatest homology (Figure 1E). For each DGR,
123 information regarding mobility, species, variable protein clustering, and domain groups was overlaid in
124 concentric rings surrounding a DGR RT phylogenetic tree (Figure 1D). As a result of their unique
125 sequence features^{10,12}, DGR RTs formed a monophyletic clade distinct from other classes of RTs,
126 consistent with prior studies^{17,19,55}. Among the 1,113 DGRs identified, 1055 (95%) resided within
127 predicted MGEs (Figure 1D, ring 1). Variable protein clustering (Figure 1D, ring 3) and domain group
128 relationships (Figure 1D, ring 4) mirrored the phylogenetic patterns of their cognate DGR RT proteins. In
129 contrast, species designations were discordant (Figure 1D, ring 2), supporting the conclusion that
130 *Bacteroides* DGRs evolve as cohesive units, encoding variable proteins that co-evolve with their
131 diversification machinery and are horizontally transferred between strains and species.

132

133 None of the *Bacteroides* variable proteins we identified had previously been annotated, therefore, we
134 used profile-based homology to infer their functions^{54,56}. VR sequences were uniformly located at the C-
135 termini of variable proteins and were predicted to adopt variant C-type lectin (C-Lec) folds as observed
136 previously^{12,19}. Nearly all variable proteins (1092/1113, 98.1%) contained structural domains with
137 homology to one of five broad groups (Figure 1E). Domain group 1 proteins have binding folds similar to
138 Type V pilins expressed by *Bacteroides*, *Porphyromonas*, and related species (PDB: 4EPS, 4QB7, and
139 5NF4)⁴⁸. This group contains a mixture of prophage-encoded (201/287, 70%), ICE-encoded (13%), and
140 cellular genome-encoded (17%) DGRs. Variable proteins within domain group 2 showed greatest

141 similarity to a *Thermus aquaticus* prophage-encoded diversified protein (TaqVP, PDB: 5VF4)⁵⁷. The
142 absence of other identifiable domains and the observation that group 2 DGRs are found within prophage
143 genomes suggests that these proteins could function as receptor binding components of phage tail fibers.
144 Domain group 3 contains large (>2,000 amino acids) multi-domain variable proteins that include motifs
145 with homology to the active regions of CotH kinases⁵⁸ (PDB: 5JDA), adjacent to other functional domains
146 such as a carbohydrate-binding motifs and leucine-rich repeats (Figure 1E). Overall, *Bacteroides* DGR
147 variable proteins display considerable modularity, whereby diversified C-terminal ligand binding
148 sequences are connected to motifs that are predicted to mediate pilus localization, association with phage
149 tail fibers, signal transduction, or other functions.

150

151 ***Bacteroides* DGRs encode pilus subunits and related variable proteins**

152 Based on their conservation and widespread distribution, we reasoned that *Bacteroides* DGRs provide
153 selective advantages to their hosts by accelerating evolution in the gut environment. To explore this, we
154 focused on a selection of five related yet non-identical DGRs present in *B. fragilis* 638R (*Bfr*), *B.*
155 *thetaiotaomicron* VPI-5482 (*Bth*), *B. uniformis* 8492 (*Bun*), *B. ovatus* 8483 (*Bov*), and *B. fingoldii*
156 *CL09T03C10* (*Bfi*) (Figure 1D arrows). Each of these DGRs diversifies variable proteins that share
157 homology to adhesive pilins located at the tips of Type V pili (domain group 1, Figure 1E). Type V pili are
158 modular, extracellular structures composed of anchor, stalk, and tip pilins (Figure 2A), with numerous
159 genes encoding homologs of each subunit type organized into operons spread throughout *Bacteroides*
160 genomes⁴⁸. For example, *Bfr* contains 93 genes clustered into 22 operons that encode components of
161 Type V pili (Supplemental Figure 2A). While the exact roles of these surface appendages in *Bacteroides*
162 spp. are at an early stage of analysis, homologous pili in *Porphyromonas* spp. facilitate coaggregation
163 with other microbes and promote colonization of the oral cavity^{59,60}.

164

165 The *Bfr*, *Bth*, and *Bun* variable protein genes (*bfrT*, *bthT*, and *bunT*) are positioned at the ends of operons
166 predicted to encode anchor and stalk proteins, which link tip pilins to the bacterial cell surface (Figure
167 2A-C, Supplemental Figure 2B). The *Bov* and *Bfi* elements are located adjacent to prophage genes, but
168 it is unclear if they are carried by a phage (Figure 1D). All five *Bacteroides* DGRs encode an Avd-like
169 protein⁶¹, and *Bfr*, *Bth*, and *Bun* encode an additional accessory factor, Msl (MuTS-like), with homology
170 to the mismatch recognition domain of MutS¹² (Figure 2B-D). Each of the five TRs contains 30 to 45
171 adenines, providing the capacity to generate massively diverse repertoires of 10¹⁸-10²⁷ potential VR DNA
172 sequences, and 10¹⁴-10²⁵ different polypeptides at the C-termini of their cognate variable proteins.

173

174 We generated predicted 3D structures with high per-residue confidence scores (pLDDT >90) for all five
175 variable proteins and compared them to known atomic structures of *Bacteroides* pili (Figure 2B-E)^{48,62}.
176 BfrT and BthT, which share 61% amino acid identity (AAI) with each other, displayed a hybrid pilin
177 structure with three domains. Their N-terminal domains (NTDs) are homologous to the NTD of BvuFim1C,
178 a structurally characterized Type V stalk pilin encoded by *B. vulgatus*⁴⁸ (Figure 2B, Supplemental Table
179 5,6). Their C-terminal domains (CTDs), which contain DGR-diversified residues organized in a C-Lec
180 fold, diverge from BvuFim1C and instead adopt a globular head structure similar to ligand-binding CTDs
181 of tip pilins (Figure 2B). Interestingly, both BfrT and BthT exhibited an additional third domain that
182 connects stalk and tip domains together. The *Bun* DGR variable protein, BunT, adopts a canonical
183 bipartite pilin structure, with N- and C-terminal domains that share high homology to BovFim1C (Figure
184 2C, Supplemental Table 5)⁴⁸, a Type V tip pilin encoded by *B. ovatus*. While the globular head of
185 BovFim1C is static, the BunT globular head displays diversifiable VR-encoded residues (Figure 2C). The
186 *Bov* and *Bfi* variable proteins, BovT and BfiT, are comparatively small, highly similar to each other (86.5%
187 AAI), and fold into a structure that is homologous to the CTD globular head of BovFim1C (Figure 2D),
188 with an N-terminal pair of alpha helices in place of the pilin-like NTD (Figure 2D, Supplemental Table 7).
189 Comparing predicted structures of *Bacteroides* VPs with structurally characterized DGR VPs^{14,15} revealed
190 a remarkable superimposition of the overall tertiary structure, including overlap in the spatial locations of
191 variable residues (Figure 2E) despite substantial differences in amino acid sequences, providing
192 evidence for the conservation of similar ligand-binding interactions¹³.

193

194 Biogenesis of a Type V pilus is a multistep process that involves: 1) pilin translocation and lipidation, 2)
195 signal peptide cleavage, 3) translocation to the outer membrane, and 4) incorporation into a growing
196 pilus^{48,63,64} (Figure 2A). As shown in Figure 2B-C, BfrT, BthT, and BunT encode conserved N-terminal
197 signal sequences, lipobox motifs, and protease-cleavable arginines required for pilus assembly. We
198 placed affinity tags⁶⁵ at the C-termini of BfrT and BfrT-C28A, a mutant derivative lacking the conserved
199 cysteine required for lipidation and translocation to the outer membrane⁶⁴, and expressed the tagged
200 proteins in *Bfr*. Following induction and cell fractionation (Supplemental Figures 3A, B, Supplemental
201 Table 8), BfrT was readily detectable in membrane fractions in both pre-processed and mature forms,
202 and in the periplasm as the mature form (Figure 2F). Treatment of intact cells with proteinase K resulted
203 in digestion of BfrT, consistent with its localization on the cell surface, while the C28A mutant was
204 protease resistant (Figure 2F). Additionally, immunofluorescence demonstrated BfrT on the cell surface
205 that was sensitive to proteinase K, but no cell surface staining of the C28A mutant was observed
206 (Supplemental Figure 3C). Mass spectroscopy of mature BfrT showed that cleavage had occurred at R43
207 (Supplemental Figure 3D, E), as predicted. We next examined the localization of tagged BovT and BfiT,

208 both of which were found exclusively in the periplasm and were resistant to proteinase K (Figure 2G).
209 These observations identify two classes of variable proteins within our five *Bacteroides* strains. DGRs
210 belonging to *Bov* and *Bfi* diversify periplasmic proteins that are structurally related to tip adhesins, but
211 either require additional factors for incorporation into pilus structures or have evolved to perform different
212 functions dependent on their ligand-binding capabilities, whereas *Bfr*, *Bth*, and *Bun* DGRs diversify Type
213 V pilus tip adhesins.

214

215 **Horizontal transfer of *Bacteroides* DGRs**

216 Horizontal transfer of DGRs that diversify phage tail fiber proteins has been well studied^{17,18}. In contrast,
217 the mobility characteristics of DGRs that target bacterial proteins are relatively unexplored. To address
218 this, we exploited the observation that the DGRs encoded by *Bfr* and *Bth* are flanked by mobility and
219 transfer genes characteristic of ICEs (Figure 3A). Like phage, ICEs often confer selective advantages to
220 their hosts by carrying cargo that encode colonization factors, metabolic capabilities, virulence
221 determinants, antibiotic resistance, or other accessory functions⁶⁶⁻⁷³.

222

223 We identified 41 DGRs residing within ICEs present in 10 different *Bacteroides* species, including the *Bfr*,
224 *Bth*, and *Bun* elements in Figure 2B,C, (Supplemental Table 4). The majority of target genes diversified
225 by these DGRs (38/41) encode variable proteins from domain group 1 that are predicted to function as
226 Type V pilus tip adhesins and are encoded directly downstream from anchor and stalk subunits
227 (Supplemental Figure 4A). Of these 41 ICEs, 16 aligned closely with each other and shared conserved
228 features including conjugation and DNA integration genes (*tra*, *int*), homologs of known transcriptional
229 regulators (*merR*, *rteC*, *araC* family members)⁷⁴, and direct repeats resulting from site-specific
230 chromosomal integration into a tRNA-Lys locus (Figure 3A, Supplemental Figure 4A). In addition to
231 related variable proteins, the DGRs encoded by these ICEs share similar TRs and RTs (Supplemental
232 Tables 9 and 10), suggesting they disseminated among *Bacteroides* species via horizontal transfer of an
233 ancestral DGR-encoding ICE.

234

235 To measure ICE activity, we developed PCR assays to differentiate integrated vs. excised forms of the
236 *Bfr* and *Bth* elements (Figure 3A). When WT cells were grown *in vitro*, only the integrated form of the
237 ICEs could be detected (Figure 3B). This was not unexpected, given that mobile genetic elements often
238 require environmental signals to induce their mobility⁷⁴. We next identified ICE-associated regulatory loci
239 (Figure 3A) and created strains that ectopically express each regulatory factor (Supplemental Table 11).
240 Overexpression of *araC2* induced ICE excision and circularization, as both circular episomes and
241 chromosomal scars were observed in *Bfr* and *Bth* (Figure 3B). Overexpression of *merR* or *rteC*, however,

242 resulted in the absence of both integrated forms and episomes, indicating the loss of ICEs from these
243 cells (Figure 3B and Supplemental Table 11). Efficient excision in the absence of integration presumably
244 leads to episome segregation during replication. To further identify requirements for excision and
245 integration, we individually deleted ICE-encoded integrase (*int1*, *int2*) and topoisomerase (*top1*, *top2*)
246 genes. Excision of the *Bfr* ICE was dependent on the presence of an intact *int2* integrase gene (Figure
247 3A, B), but was unaffected by knocking out *int1*, *top1*, or *top2* (Supplemental Table 12). Finally, to
248 determine if inducing signals are provided *in vivo*, we colonized germ-free Swiss Webster mice with *Bfr*
249 and measured relative levels of episome and chromosomal scar formation by qPCR in bacteria recovered
250 from fecal pellets, cecal contents, and colonic mucosa. As shown in Figure 3C, compared to *Bfr* plus
251 empty vector cultured *in vitro*, we observed 17- and 27-fold increases in ICE activity in fecal pellets and
252 cecal content, respectively, and a nearly 700-fold increase in colonic mucus, the natural habitat of
253 *Bacteroides*.

254
255 Resident ICEs are known to exclude integration by homologous mobile elements^{75,76}. To create recipient
256 cells suitable for mating experiments, we overexpressed *merR* to promote ICE excision and loss (Figure
257 3B), curing *Bfr* and *Bth* strains of their DGR-containing ICEs (Δ ICE) and leaving their chromosomal
258 integration sites free. In matings between isogenic WT donor and Δ ICE recipients, transconjugants were
259 isolated at a frequency of 10^{-5} to 10^{-7} , with *Bth* donors displaying greater transfer efficiencies than *Bfr*
260 donors (Figure 3D, E, Supplemental Table 13). To determine if transfer requires the DGR-encoded
261 variable protein, we deleted *bfrT* and found that *Bfr* Δ *bfrT* mutants displayed the same transfer efficiencies
262 as the WT parent, demonstrating that the diversified Type V pilus tip adhesin is dispensable for horizontal
263 transfer. ICE conjugation was also observed in gnotobiotic mice (Supplemental Figure 4B,C). Our results
264 demonstrate that DGRs in *Bfr* and *Bth* are encoded within functional ICEs that undergo conjugative
265 transfer *in vitro* and *in vivo*, providing an explanation for the phylogenetic distribution of DGRs in
266 *Bacteroides* and a mechanism for the horizontal transfer of accelerated protein evolution between
267 species.

268

269 **Differential control and mechanistic conservation of mutagenic retrohoming**

270 To characterize the real-time dynamics and mutational patterns of *Bacteroides* DGRs, we measured
271 mutagenic retrohoming levels *in vitro* and *in vivo* and interrogated the diversified sequences. Strains
272 carrying the DGRs shown in Figure 2B-D were grown *in vitro*, sampled over a two week period, and VRs
273 were barcoded, amplified, and deep sequenced (Figure 4A). We calculated the percentage of VRs that
274 had diverged from their parental sequence and found that the *Bov* and *Bfi* elements showed remarkably
275 high levels of mutagenesis, with 13% or 40% of VRs diversified by day 14, respectively (Figure 4B). This

276 was unexpected, since the activity of DGRs that mutagenize bacterial genes in other genera has been
277 reported to be low or absent during *in vitro* growth, reflecting their apparent regulation^{10,19}. VR
278 mutagenesis was abolished in strains harboring knockout mutations in *rt* (Δrt) (Figure 4C, Supplemental
279 Figure 5A) and restored by complementation with wild type *rt* expressed at an ectopic location
280 (Supplemental Figure 5B). In contrast, VR sequences from *Bfr*, *Bth*, and *Bun* displayed *in vitro* levels of
281 mutagenesis that were 100- to 10,000-fold lower than *Bov* or *Bfi* (Figure 4B). To measure DGR activity
282 *in vivo*, we monocolonized germ-free Swiss Webster mice with individual *Bacteroides* strains.
283 Colonization levels in the GI tract were similar for each *Bacteroides* strain as measured by colony-forming
284 units in fecal pellets (Supplemental Figure 5C). *Bov* and *Bfi* DGRs were highly active in the murine GI
285 tract, displaying levels of diversity similar to those observed *in vitro*, while the activity levels observed
286 with *Bfr*, *Bth*, and *Bun* remained low (Figure 4B). Next, we used RNA-Seq to probe relationships between
287 mutagenic retrohoming and transcription of DGR-encoded genes. We observed significantly higher
288 relative amounts of transcripts encoding *avd*, TR, and *rt* in high activity strains (*Bov* and *Bfi*, Figure 4D)
289 compared to those with low DGR activity (*Bfr*, *Bth*, *Bun*), suggesting that DGR mutagenesis is regulated,
290 at least in part, at the transcriptional level. Thus, mutagenic retrohoming levels in our five *Bacteroides*
291 strains fall into two categories. DGRs carried by *Bov* and *Bfi* are constitutively active, while those in *Bfr*,
292 *Bth*, and *Bun* appear to be tightly regulated, as commonly observed in other systems^{19,77}.

293
294 Mutagenic retrohoming in *Bacteroides* demonstrated remarkable specificity for substitutions at TR
295 adenines (Figure 4E), a hallmark of DGR RT enzymes observed across taxa^{16,18,77}. To identify positional
296 effects, we exploited the constitutive activity of the *Bov* and *Bfi* DGRs to examine time-dependent levels
297 of mutagenesis as a function of position within VR (Figure 4F, Supplemental Figure 5D). The central
298 region of VR is enriched with sequences corresponding to TR AAC motifs (Figure 4F, underlines) that
299 enable random substitution at one or both adenines^{12,20}, accounting for the significant accumulation of
300 mutations over time. In contrast, VR positions at the 5' and 3' ends displayed minimal to no mutations
301 despite the presence of TR adenines, consistent with boundary effects similar to those reported for the
302 *Bordetella* phage BPP-1 DGR¹⁶, including a crossover interval near the center of the 3' boundary region
303 where mutagenized cDNA integrates to replace parental VRs. Finally, by counting the number of adenine
304 mutations in uniquely diversified VR sequences, we observed a median mutational density of about 50%
305 of available positions (Figure 4G, Supplemental Figure 5F). A nearly identical mutational density was
306 reported in *Bordetella in vivo*⁷⁸, and *in vitro* with cDNA synthesized by purified BPP-1 RT, Avd, and TR-
307 RNA in the presence of dNTPs^{79,80}. These observations highlight the striking conservation of the
308 mechanisms of adenine-specific mutagenesis and cDNA integration in distantly related DGRs and
309 bacterial hosts.

310

311 ***Bacteroides* DGRs preferentially create non-synonymous substitutions that alter side chain**
312 **chemistry**

313 The high levels of mutagenic retrohoming conferred by the *Bov* and *Bfi* DGRs allowed us to build the
314 largest dataset of experimentally derived diversified VR sequences available to date. VR-encoded amino
315 acids that differed from their initial parental sequence during serial subculturing were classified as
316 synonymous or non-synonymous, and the amino acid substitution frequency at each codon position was
317 calculated. Non-synonymous substitutions at cognate TR adenine positions predominated our dataset,
318 accounting for 99.9% of amino acid changes resulting from nearly 4,500,000 DGR-generated VR
319 mutations (Figure 5A, Supplemental Figure 6A-C). The vastly disproportionate number of non-
320 synonymous substitutions generated by DGR mutagenesis arises directly from the bacterial genetic code,
321 coupled with a high abundance of TR AAY (AAC or AAT)¹⁵ motifs, which account for 10 of the 12 variable
322 codons between the 5' and 3' boundaries of the *Bov* VR (Figure 4F). The abundance of TR AAY motifs
323 extends beyond *Bov*. In our *Bacteroides* dataset, AAY motifs (~10.5/TR) outnumber single-adenine
324 motifs (~5.1/TR) (Supplemental Figure 6D,E Mann-Whitney $p < 0.0001$), underscoring their significance.
325 For AAC motifs, 16 potential codons can be generated through random adenine mutagenesis, 14 of which
326 encode unique amino acids (Figure 5B), while only two codons are synonymous (AGC and TCC which
327 encode serine), and a stop codon can never be generated. The side chains of the 15 amino acids
328 produced by AAC mutagenesis encompass the entire range of available chemical properties, including
329 polar uncharged, small hydrophobic, large hydrophobic, positive charge, negative charge, and no side
330 chain (glycine) (Figure 5B). Therefore, the predominance of TR AAY motifs leads to an expansive and
331 chemically diverse list of amino acids at variable sites and an overwhelmingly high probability of
332 producing non-synonymous mutations due to adenine mutagenesis.

333

334 On closer examination, we noticed a non-random pattern of diversified residues whereby the chemical
335 property of the side chain often switched during mutagenesis (Supplemental Figure 6F). By positioning
336 an adenine in the middle of the codon, random mutagenesis of the first position results in four potential
337 codons, each with unique side chain chemical properties (Figure 5B, arrow). Interestingly, we observed
338 a non-random pattern in the nucleotide frequency at variable positions. Adenine was the most prevalent
339 nucleotide at the majority of the variable positions in both *Bov* (Figure 5C, D) and *Bfi* (Supplemental
340 Figure 6G, H), followed by guanine and thymine, while cytosine was rarely observed. When specifically
341 focusing on VR positions corresponding to TR AAC motifs, an adenine was observed in the second
342 position at very high frequencies that were often greater than 50%. Thus, the nucleotide frequency across

343 variable positions in VR exhibits stochastic but non-uniform patterns, with a clear bias to incorporate
344 adenines that result in frequent switching of amino acid side chain chemistry.

345

346 Mechanistic models of mutagenic retrohoming predict the formation of heteroduplexes between
347 mutagenized cDNAs and parental VR sequences with an unusually high density of mismatches^{10,78,80}.

348 Taking advantage of the constitutive activity of the *Bov* DGR, we explored the impact of mismatch repair
349 on the level and pattern of VR mutagenesis by knocking out *mutS*. *Bov* Δ *mutS* exhibited substantially
350 higher levels of VR mutagenesis than WT, with up to 25% and 43% of VRs mutated by days 3 and 14,
351 respectively, corresponding to an approximate 5- to 10-fold increase compared to the parent strain
352 (Figure 5E). Complementation with intact *mutS* restored mutagenic retrohoming to WT levels
353 (Supplemental Figure 6I). Intriguingly, analysis of nucleotide frequencies at mutated variable sites
354 revealed an almost identical distribution (Figure 4G) and pattern (Figure 5C, F) when diversified VRs
355 were compared between *Bov* WT and *Bov* Δ *mutS*. These results support a mechanism in which MutS-
356 mediated repair operates in an all-or-nothing manner, whereby VR heteroduplexes are either converted
357 back to the parental sequence or fully escape mismatch repair.

358

359 **DGR dynamics under competitive pressure**

360 Estimating the percentage of VRs that have undergone mutagenic retrohoming provides an incomplete
361 picture of the true extent of sequence diversity. For example, a population in which a majority of diverged
362 VRs encode the same or a limited number of DNA or protein sequences would clearly be distinct from
363 one in which most of the diverged sequences differed from each other. Thus, we calculated the Shannon
364 entropy⁸¹ of VRs that had undergone adenine-mutagenesis under different conditions. Although this
365 metric can be biased at small sample sizes⁸², our large mutational dataset allows it to encompass the
366 relative strengths of two opposing forces: i) mutagenic retrohoming (Supplemental Figure 6J), which
367 increases VR entropy by randomizing sequences, and ii) purifying selection, which decreases entropy
368 through preferential propagation of mutagenized VR sequences that provide a competitive advantage¹⁹.

369

370 We measured VR entropy in populations of *Bov* that were grown *in vitro* or in germ-free mice colonized
371 with or without Altered Schaedler Flora (ASF), an eight-member bacterial consortium, to provide
372 interspecies competition⁸³ (Figure 5G). As expected, VR sequences derived from *in vitro* grown cells
373 displayed high levels of amino acid entropy at early and late timepoints (Figure 5H), indicating that
374 mutagenic retrohoming was primarily driving VR diversity. When *Bov* was introduced into germ-free mice,
375 VR entropy displayed high values similar to *in vitro* samples, suggesting the absence of strong selective
376 forces in monocolonized hosts. In contrast, VR sequences from animals co-colonized with both *Bov* and

377 ASF displayed a time-dependent decrease in entropy that became highly significant by week 2 post-
378 gavage, indicative of positive selection. On closer examination, samples from *in vitro* grown or
379 monoassociated *Bov* contained VRs that were almost entirely different from each other, while samples
380 from mice co-colonized with ASF contained populations in which two to three unique VRs comprised
381 >50% of all mutated sequences (Figure 5I). Furthermore, in separately caged co-colonized mice, mutated
382 populations of VRs had converged by day 14 to express similar amino acid sequences despite being
383 encoded by different DNA sequences. In the examples shown in Figure 5J, several hydrophobic side
384 chains have been replaced by polar, uncharged residues while maintaining a hydrophobic interaction site
385 near the top of the VR, creating a more open binding pocket and suggesting these variant sequences
386 have evolved to interact with a new, common ligand. Taken together, these observations demonstrate
387 that in the face of competition with other microbes, VR entropy decreases in a manner expected for
388 environmental conditions that exert positive selection.

389

390 **DGRs are active during the intergenerational handoff from mothers to infants**

391 We hypothesized that dynamic changes that accompany the intergenerational handoff of gastrointestinal
392 microbes from mothers to infants during birth could select for DGR-mediated adaptations that facilitate
393 colonization and persistence in a new host. To explore the role of DGRs during this period, we performed
394 an integrative analysis of retroelements present in metagenomic datasets of human fecal microbiomes
395 from 144 longitudinally sampled mother-infant pairs^{84–86} as well as from a dataset of 146 healthy adults⁸⁷
396 (Supplemental Figure 7A). We identified 5106 different DGRs, of which 2740 were identified within
397 mother-infant pairs, with 698 DGRs found uniquely in infants, 1654 found only in mothers, and 388 DGRs
398 that were apparently transmitted from mothers to infants during or after birth (Figure 6A, Supplemental
399 Table S14). To gain a deeper understanding of these elements, we analyzed DGR phylogeny, variable
400 protein domains, and predicted transfer vector, similar to our prior analysis with *Bacteroides* (Figure 1)
401 but expanded to include all phyla (Figure 6B).

402

403 The majority of DGRs in this dataset were identified as phage- or prophage-encoded (93%). Of the 344
404 DGRs found in cellular genomes, plasmids, or ICEs, 214 (62%) were predicted to reside within
405 Bacteroidota, while 106 (31%) were classified as Bacillota. Cellular or ICE-encoded Bacteroidota DGRs
406 almost exclusively diversify pilus subunits or other cell adhesion proteins, as observed with the dataset
407 analyzed in Figure 1B. In contrast, most Bacillota variable proteins classified as cellular, plasmid, or ICE-
408 encoded have domains similar to phage receptor binding proteins, suggesting this binding module was
409 co-opted to perform some other function. DGRs in this dataset were distributed throughout infants,

410 mothers, and adults, except for a group of DGRs belonging to Actinomycetota (Figure 6B, red star), which
411 were enriched in infants, and to a lesser extent in mothers, but rarely found in nonpregnant adults.

412

413 Mode of delivery had a significant impact on the number of DGRs identified in infants (Figure 6C). At
414 birth, DGRs were much more common in vaginally born infants compared to those born by C-section,
415 and this trend persisted throughout the first year of life. The mean number of unique DGRs in infants at
416 one year was significantly less than the mean number of DGRs in mothers and adults (Figure 6C,D,
417 $p < 0.0001$ all comparisons of 1-year-olds vs adults), showing that new DGRs continue to be acquired
418 throughout life. Mode of delivery also correlated with the relative distribution of DGR-containing taxa.
419 Prior to delivery, there were no significant differences in the taxonomy of microbes harboring DGRs in
420 mothers undergoing C-section compared to mothers delivering vaginally (Figure 6E). However, infants
421 delivered vaginally showed a much higher proportion of DGR-containing Bacteroidota, while infants born
422 via C-section acquired an initial set of DGRs that were almost exclusively encoded by Bacillota. By the
423 end of infants' first year, the taxonomy of DGR-containing microbes more closely resembles the
424 distribution found in adults regardless of the mode of delivery (Figure 6E). Comparisons between
425 breastfed and formula fed infants or between males and females revealed no differences in the number
426 of DGRs or their taxonomic distribution (Supplemental Figure 7B-E). Together, these results demonstrate
427 that mode of delivery has a significant impact on the number and types of microbes harboring DGRs,
428 with vaginally born infants acquiring a larger number that more closely resemble the phylogenetic
429 distribution observed in adults.

430

431 To determine if DGRs were active at any point within these samples, we aligned raw sequencing reads
432 with identified VRs and searched for adenine-specific mutations. The percentage of active DGRs in
433 infants varied from 56% to 80% over the first year of life but was similar to maternal and adult levels
434 (Figure 6F). Next, we calculated the consensus VR amino acid sequence at each timepoint and compared
435 VR haplotypes between mothers and infants. Of 388 transmitted DGRs, 72% showed evidence of VR
436 haplotype switching (Figure 6G), raising the possibility that new variable proteins had been selected in
437 the infant. A representative example of a diversified type V pilin homolog that was present in a mother at
438 birth, and subsequently detected in her infant's profile over the first year of life, is illustrated in Figure 6H.
439 Comparing VR sequences shows that a majority of variable codons underwent continued alterations that
440 changed the chemical class of diversified residues at four months, stabilizing by 1 year.

441

442 DGRs that were transmitted from mothers to infants born vaginally were most likely to be found in
443 Bacteroidota (56%), rather than Bacillota (15%) (Figure 6E), and they were more likely to diversify pilus

444 proteins than the general population of DGRs (15% vs. 10%, Chi-Square $p < 0.0001$). Importantly, we
445 were able to identify DGRs in our dataset that were nearly identical to each of the *Bacteroides* elements
446 depicted in Figure 2B-D. Although the *Bfr*, *Bth*, and *Bun* DGRs we examined were quiescent under
447 laboratory conditions (Figure 4B), homologous elements were highly active in mothers and infants,
448 showing clear evidence of adenine-templated VR mutagenesis (Supplemental Figure 7F-G,
449 Supplemental Table S15). For example, Figure 6I shows TR and VR sequences from a *Bun* DGR
450 homolog that diversifies a Type V pilus tip adhesin, and the predominant VR haplotypes generated by a
451 DGR that was transferred from mother to infant. As expected, elements homologous to the *Bov* and *Bfi*
452 DGRs that were highly active *in vitro* and in germ-free mice were similarly active in humans.

453

454 These observations begin to characterize the abundance, distribution, and activities of DGRs in human
455 infant and adult populations and their ability to diversify a wide array of potential ligand binding proteins,
456 including pilus-associated adhesins. We also provide evidence that DGRs are transferred during the
457 intergenerational handoff, that mode of delivery profoundly influences the relative abundance of DGR-
458 containing microbes in the newborn gut, and that DGR transfer is associated with the appearance of new
459 VR haplotypes that predominate in the infant's gastrointestinal tract following maternal transmission, as
460 would be expected for genotypic alterations that are subject to positive selection.

461

462 Discussion

463 Numerous studies have highlighted the impact of microbial genetic variation on human health and
464 disease^{88–90}. This variation is driven by rapid adaptations that can quickly spread throughout microbial
465 communities^{91–93}. DGRs play a unique role in variation and adaptation by creating hypermutable hotspots
466 that produce unrivaled levels of protein diversity, through a mechanism that targets ligand-binding
467 residues and is shared across vast phylogenetic distances. By characterizing the variable proteins,
468 mutational dynamics, and modes of transmission of *Bacteroides* DGRs, we can build a foundation for
469 understanding how genome plasticity shapes host-associated microbial communities.

470

471 Of over 1,100 unique, DGR encoded variable proteins identified in human gut-associated *Bacteroides*,
472 nearly 25% are predicted to diversify Type V pilins^{48,63}. The genetic systems responsible for these
473 structures display a remarkable degree of modularity and apparent redundancy, with the number,
474 location, and organization of genes encoding pilin homologs differing between human-associated species
475 and strains^{48,63}. Less than half of the 16 multigene operons that encode pilin subunits in *Bfr* show
476 evidence of protein expression *in vitro* (Supplemental Figure 2A), and many of these gene clusters are
477 flanked by integrases, transposases, prophage loci, transfer genes and other signs of mobility⁴⁸. Thus,
478 differential expression and horizontal transfer may partly explain the modular complexity observed in
479 *Bacteroides* genomes. For DGRs that target adhesive pilins in *Bfr*, *Bth*, and related taxa with homologous
480 ICEs, we propose that two factors promote their dissemination and positive selection. The first involves
481 properties of the conjugative elements, such as inducibility in mucus and the presence of stalk and anchor
482 homologs that could help display diversified pilin tips. The second is the availability of new hosts with
483 genetic backgrounds that are diverse and adapted to acquire horizontally transferred type V pilus genes.
484 Adhesive pili are often observed to be essential determinants of microbe-host interactions involved in
485 colonization, and microbe-microbe interactions that structure bacterial communities and facilitate biofilm
486 formation^{59,60}. An understanding of the selective advantages conferred by type V pili will benefit from the
487 construction of isogenic mutants that are completely devoid of surface pili or lack specific components,
488 the availability of animal models that recapitulate adhesive interactions in the human GI tract, and the
489 identification of ligands recognized by static as well as DGR-diversified tip adhesins.

490

491 While the abundance of DGRs in nature attests to their selective advantages^{12,19,55}, it is unknown how
492 host cells balance the benefits of accelerated evolution with the increased potential for loss of fitness. It
493 is reasonable to expect that mutagenic retrohoming will often be subject to regulation, with bursts of
494 mutagenesis strategically deployed during times of stress, population expansion, nutritional changes or
495 other factors, and interspersed with periods of quiescence that allow selection and fixation of adaptive

496 traits¹⁹. Considering this, it was not unexpected to observe minimal activity with the *Bfr*, *Bth* and *Bun*
497 DGRs *in vitro* or in gnotobiotic mice. In contrast, the constitutively high levels of activity measured with
498 the *Bov* and *Bfi* elements was surprising, and provided an opportunity to examine DGR mutagenesis at
499 a level of resolution that had not been previously attained. We observed an overwhelming (>1,000-fold)
500 preference for introducing nonsynonymous substitutions in VR, a misincorporation bias that favors
501 changes in the chemical properties of side chains available for ligand binding, and a role for host-encoded
502 MutS that supports an all-or-nothing model for repair of heteroduplex intermediates during mutagenic
503 retrohoming. In GF mice, we measured a decrease in the Shannon entropy of *Bov* VR sequences that
504 was dependent on the presence of competing microbes and resulted in the appearance of new
505 predominant VR haplotypes, as expected for conditions that favor positive selective sweeps.

506

507 Although signals and regulatory mechanisms that control mutagenic retrohoming await discovery, close
508 homologs of all of the DGRs we studied were identified in human metagenomes and observed to be
509 active in the human gut. This highlights the importance of examining these elements in their natural
510 context. For DGRs predicted to diversify bacterial factors, type V pilus subunits were the most common
511 variable proteins encoded by Bacteroidota, and they were the most likely DGRs to show evidence of
512 maternal-infant transfer. Vaginally born infants had a greater number of DGRs compared to infants born
513 by C-section, and they were mainly encoded by *Bacteroides* species as opposed to Bacillota. This
514 difference resolved quickly and mirrors well-characterized time-dependent effects of birth-mode on the
515 overall composition of the developing infant microbiome⁹⁴. Most interesting, however, is that for the
516 majority of variable proteins identified as being transferred, the predominant VR haplotypes observed in
517 mothers switched to new predominant haplotypes in their infants, consistent with the hypothesis that
518 DGR-driven adaptations are occurring during the first year of life.

519

520 Hypervariable systems and assessments of purifying selection may provide a means for identifying genes
521 and networks in microbial communities that confer contextually significant fitness advantages.
522 Accordingly, our working hypothesis is that DGRs that diversify bacterial proteins function, at least in part,
523 to optimize colonization factors that promote engraftment and maintenance of host bacteria within the GI
524 microbiota. If true, understanding DGRs and the variable proteins they diversify will not only have
525 applications for understanding microbe-host interactions, but may also provide a means to engineer
526 therapeutic microbial consortia capable of rapidly evolving colonization factors to promote efficient
527 engraftment in new hosts. This could pave the way for future applications that harness the adaptive
528 properties of DGRs to support health and reverse microbiome-associated diseases.

529

530 **Limitations**

531 There are two material limitations of our study that reflect broader challenges in efforts to understand
532 diversity in human-associated microbial communities. The first involves the need to identify phenotypes
533 from genotypic information, a transformative resource available on a massive scale that is rarely sufficient
534 to provide causal links. DGRs were discovered based on a phenotype, tropism switching by *Bordetella*
535 phage, and the genetic basis was characterized with relative ease^{17,18}. In contrast, identifying the function
536 of an uncharacterized gene, or the advantage of diversifying it, may require recapitulating natural
537 environments that provide appropriate selective pressures, including signals for expression, receptors for
538 ligand interactions, and many other context-dependent parameters. The 'genotype to phenotype to
539 mechanism' pathway can pose major challenges, but they are not necessarily insurmountable. In the
540 case of DGRs and similar systems, for example, hypervariability can provide sensitive, real-time
541 estimates of positive selection and the environmental conditions under which it occurs – information that
542 can be used to target hypothesis-driven discovery.

543

544 A second limitation involves the nature of metagenomic data available from existing large-scale efforts to
545 characterize human microbiomes, which is primarily derived from short read approaches that make
546 assembly difficult and rarely reach sufficient depth to fully capture VR sequence variability and entropy.
547 Future efforts that incorporate long-read sequencing of sufficient depth on well curated longitudinally
548 obtained mother-infant samples, along with Amplicon-Seq to deeply characterize diversified VRs, will be
549 required to understand DGR dynamics following birth and to identify variable proteins subject to positive
550 selection in developing infants.

551

552 **Acknowledgements:**

553 We thank Elaine Hsiao and her laboratory members Kristie Yu and Jorge Paramo at UCLA, as well as
554 Sarkis Mazmanian at the California Institute of Technology for valuable advice and assistance with
555 gnotobiotic mouse husbandry. We also thank Eric Martens at the University of Michigan for providing
556 several of the *Bacteroides* strains used in this study. We are indebted to Suzanne Devkota at Cedars
557 Sinai Medical Center, Partho Ghosh at the University of California, San Diego, and Blair Paul at the
558 Marine Biological Laboratory, Woods Hole, for their insightful critiques of the manuscript. This research
559 was supported by the NIH (1K08DK138316/5K12HD111040; 5K12HD000850 to B.R.M.) and the JDS
560 Family Foundation and Kavli Endowment (to J.F.M.). We acknowledge the use of resources at the UCLA
561 Proteomics Laboratory for mass spectroscopy and the UCLA Neuroscience Genomics Core for RNA-
562 sequencing. The funders had no role in the design of the study, in the collection, analyses, or
563 interpretation of data, in the writing of the manuscript, or in the decision to publish the results.

564

565 **Author contributions:**

566 Conceptualization, B.R.M., Y.W., U.A., and J.F.M, with input from all authors; Methodology, B.R.M, Y.W.,
567 U.A., and J.F.M; Investigation, B.R.M., Y.W., C.W., J.R., K.S., C.B., S.M., and U.A.; Software, B.R.M.
568 and B.M.S.; Writing – Original Draft, B.R.M. and J.F.M.; Writing – Review and Editing, B.R.M., U.A., and
569 J.F.M.; Funding acquisition, B.R.M. and J.F.M;

570

571 **Declaration of interests:**

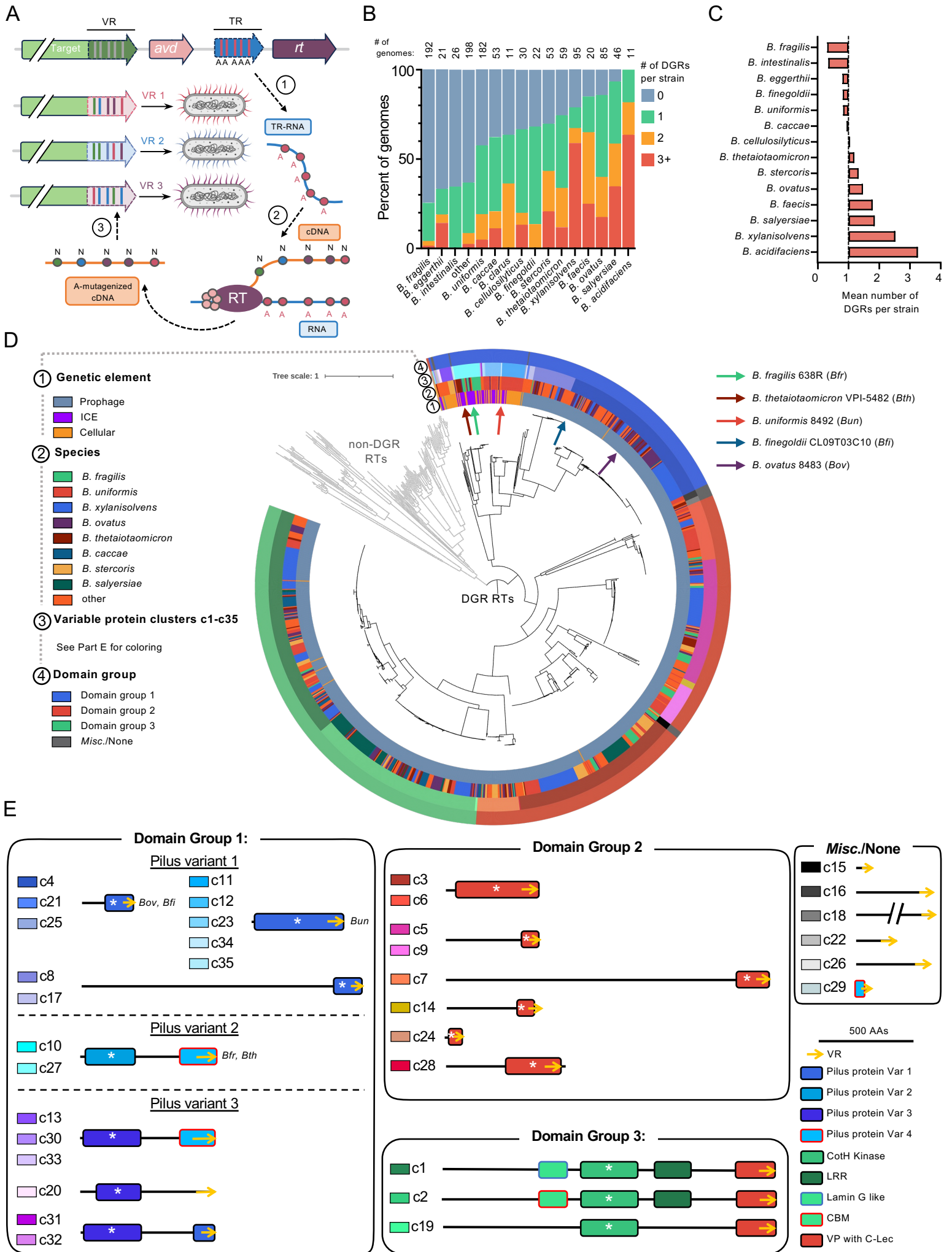
572 J.F.M. is a cofounder and chair of the scientific advisory board (SAB) of Pylum Biosciences, Inc., a
573 member of the SAB of Notitia Biotechnologies, and an advisory board member of Seed Health, INC.

574

575 **Supplemental information:**

576 Figures S1–S7.

577 Excel file containing Tables S1-S15.



578 **Figure 1. DGRs are widely distributed in *Bacteroides***

579 (A) Schematic representation of the mechanism of DGR-mediated mutagenic retrohoming. A typical
580 DGR locus is depicted on top. The target gene (green), which encodes a variable protein, contains
581 a VR sequence which is the recipient of mutagenesis. An accessory protein, Avd (pink), helps
582 guide the DGR-encoded RT (purple) to the TR-RNA template. Steps in mutagenic retrohoming
583 include: 1) production of TR-RNA; 2) error-prone cDNA synthesis with misincorporation at TR
584 adenines; and 3) cDNA replacement of parental VR alleles^{10,16,17,95}.

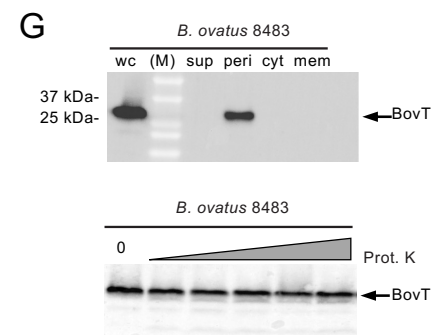
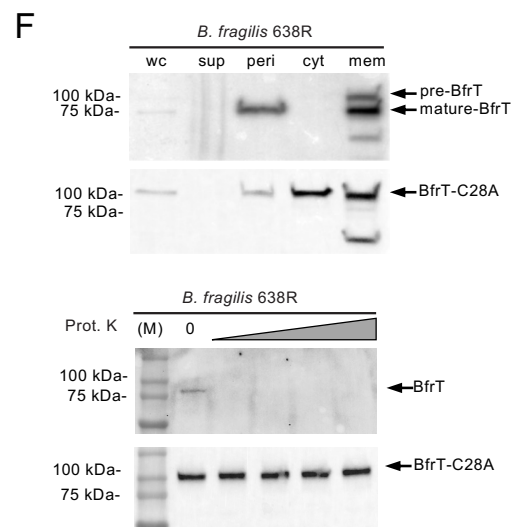
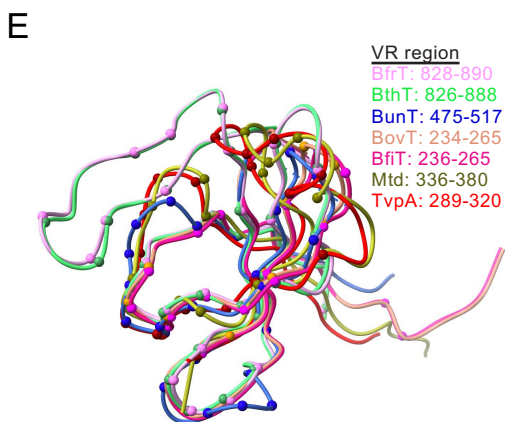
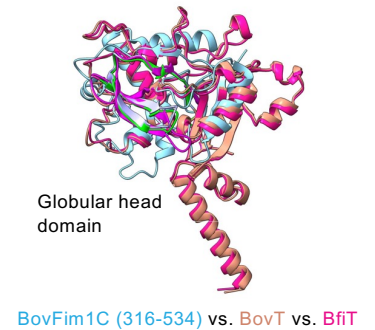
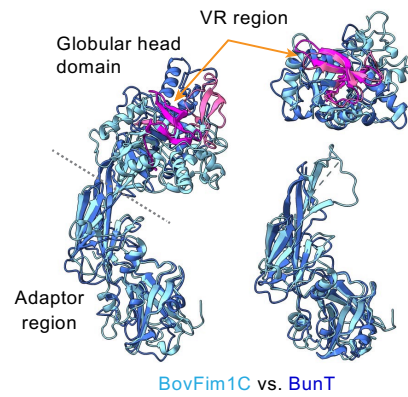
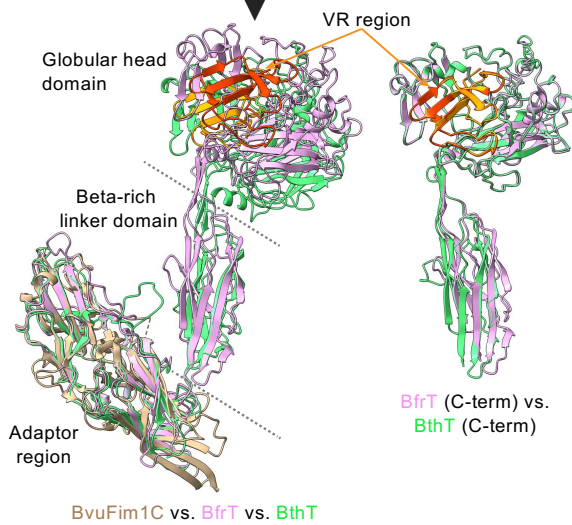
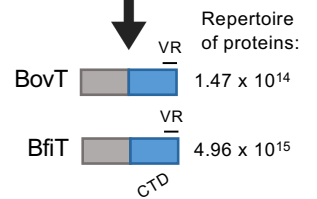
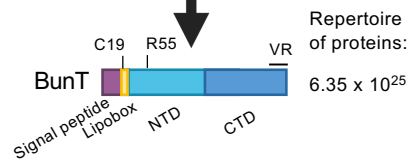
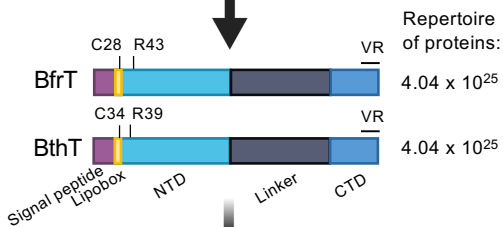
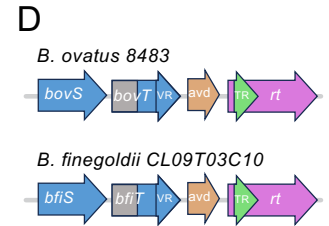
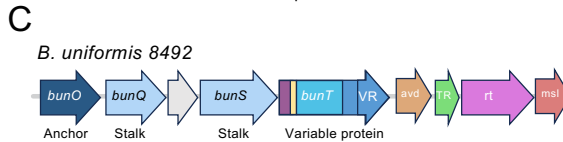
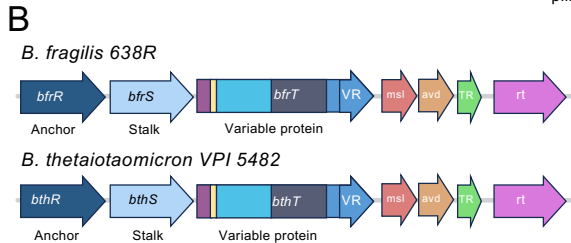
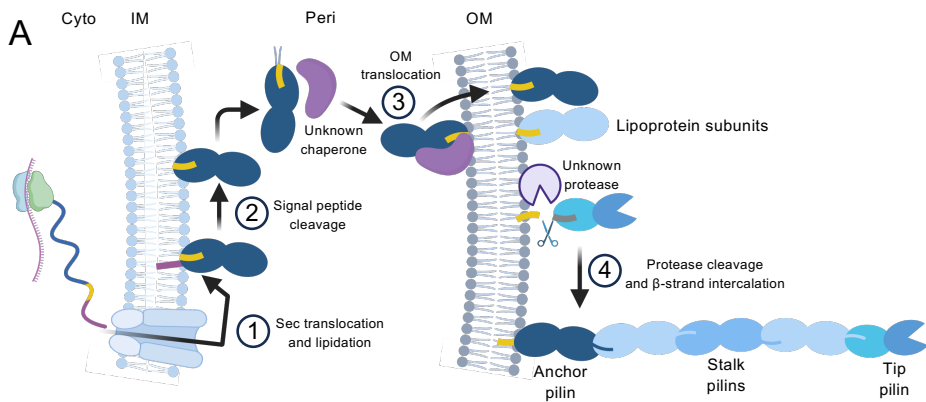
585 (B) The number of DGRs found within *Bacteroides* strains, grouped by species.

586 (C) The mean number DGRs per strain grouped by species. The average number of DGRs in all
587 strains (1.01) is set as the baseline.

588 (D) Phylogeny of DGR RTs and non-DGR RTs in *Bacteroides* genomes. Rings depict the genomic
589 location of the DGR (ring 1), species classification (ring 2), variable protein cluster (ring 3), and
590 variable protein domain group (ring 4). Variable protein clusters (c1-c35) and domain groups are
591 colored according to Part E.

592 (E) Visual representation of each cluster, grouped by domain group. The most significant domain is
593 denoted by the white asterix. The color of each cluster corresponds to Part D, ring 3.

594 See also Supplemental Figure S1 and Supplemental Table S4.



595 **Figure 2. *Bacteroides* DGRs diversify pilus tip adhesins and related proteins**

596 (A) Cartoon representation of the steps in Type V pilus assembly⁴⁸. (1) N-terminal signal peptide
597 recognition and pilus subunit translocation across the inner membrane via the SecYEG
598 translocon, (2) lipidation at a conserved cysteine residue within a lipobox motif and signal peptide
599 cleavage, (3) translocation of the lipidated protein across the outer membrane via an LPP
600 transporter, and (4) incorporation into a growing pilus structure through protease-assisted
601 cleavage which simultaneously releases the protein from the membrane and creates an acceptor
602 site for beta-strand intercalation from an incoming subunit^{48,63}. Cyto: cytoplasm; IM: inner
603 membrane; Peri: periplasm; OM: outer membrane; NTD: N-terminal domain; CTD: C-terminal
604 domain.

605 (B-D) DGR loci and their upstream genes in *Bfr*, *Bth*, *Bun*, *Bov*, and *Bfi* isolates chosen for this study.
606 Below each loci is a graphical representation of the variable proteins where each domain is
607 colored and labeled with the predicted function. The number of variable proteins that can be
608 generated through mutagenic retrohoming is specified to the right. The protein structures at the
609 bottom are the superposition of the AlphaFold⁶² structures of the variable proteins with known
610 pilus protein structures. BvuFim1C (PDB: 4QB7⁴⁸), BovFim1C (PDB: 4EPS⁴⁸).

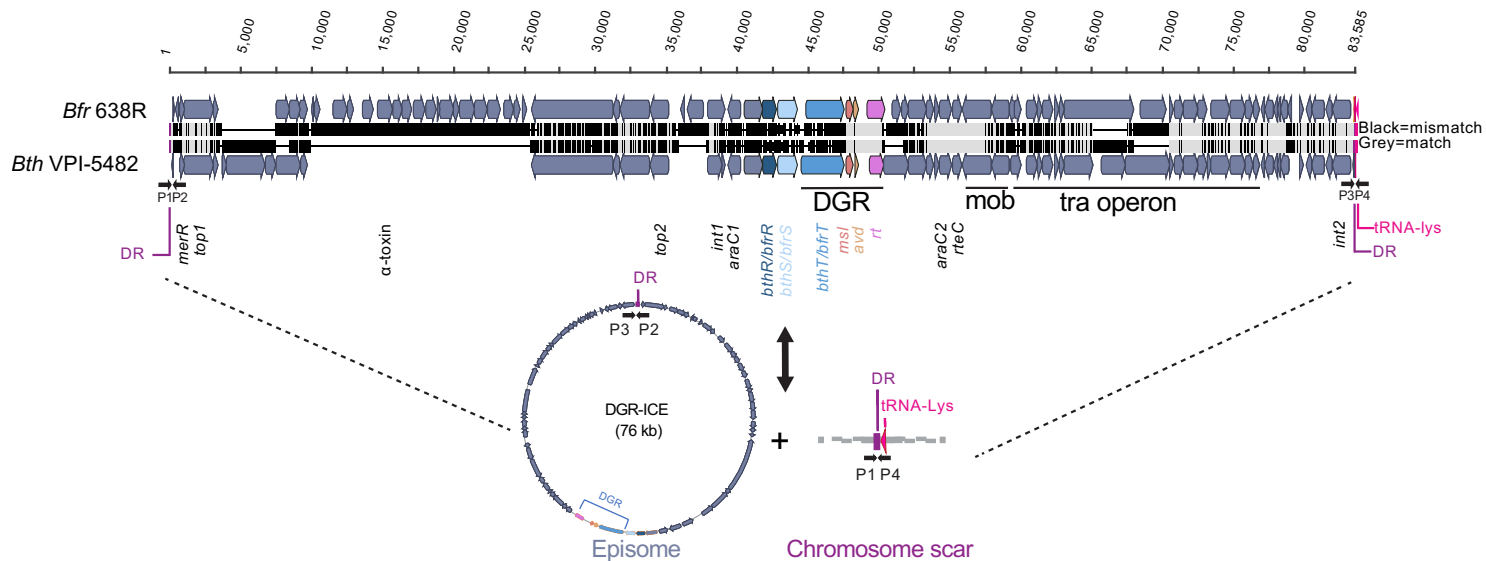
611 (E) Superposition of predicted VR-encoded structures for the five *Bacteroides* variable proteins in
612 parts (B)-(D), with the VR-encoded structures of Mtd (gold) and TvpA (red)¹³⁻¹⁵. The positions of
613 variable residues, which are often superimposable in space, are shown as colored balls.

614 (F) Immunoblot of BfrT and BfrT-C28A overexpressed in *Bfr* after cellular fractionation (top) or after
615 intact whole cells were exposed to proteinase K (bottom). See STAR Methods for fractionation
616 protocol and Supplemental Figure 3A for fractionation controls.

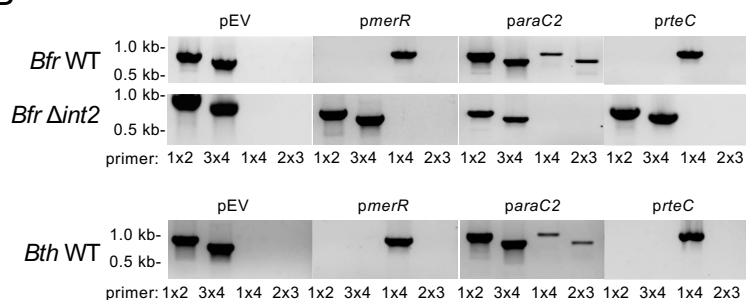
617 (G) Immunoblot of BovT overexpressed in *Bov* after cellular fractionation (top) or after whole cells
618 were exposed to proteinase K (bottom).

619 See also Supplemental Figures S2 and S3 and Supplemental Tables S5-8.

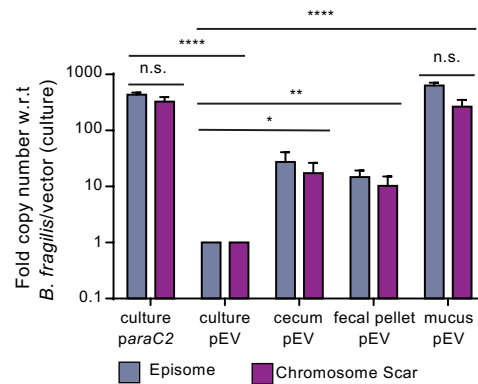
A



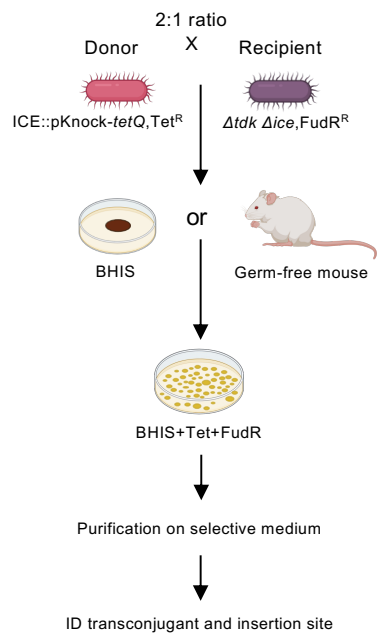
B



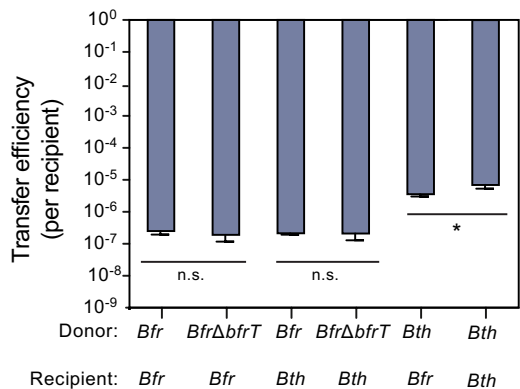
C



D



E



620 **Figure 3. *Bacteroides* DGRs are horizontally transferred between strains and species**

621 (A) (Top) Synteny between *Bfr* and *Bth* ICEs. Important elements within each ICE, including the DGR
622 loci, are labeled. (Bottom) Cartoon schematic of the *Bfr* ICE conversion between chromosomally
623 integrated and episomal forms. Binding of primers P1-P4 shown for integrated ICE, episomal ICE
624 and chromosomal scar.

625 (B) PCR products from *Bfr* (top) and *Bth* (bottom) using primers to differentiate chromosomally
626 integrated ICEs from excision products (episome and scar) following overexpression of
627 designated ICE encoded regulatory genes. pEV, empty vector.

628 (C) Ratios of excised episomes or chromosomal scars to integrated ICEs, normalized to pEV
629 containing strains cultivated *in vitro*. Samples from cecum, fecal pellet, and mucus originate from
630 monocolonized SW mice (n=3).

631 (D) Experimental setup for ICE mating assays.

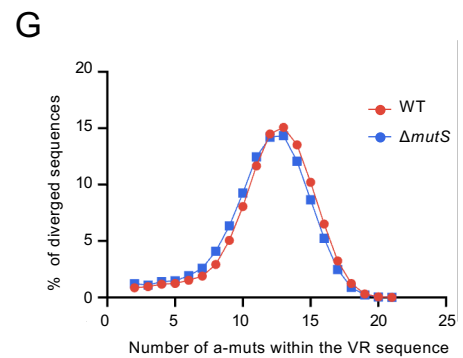
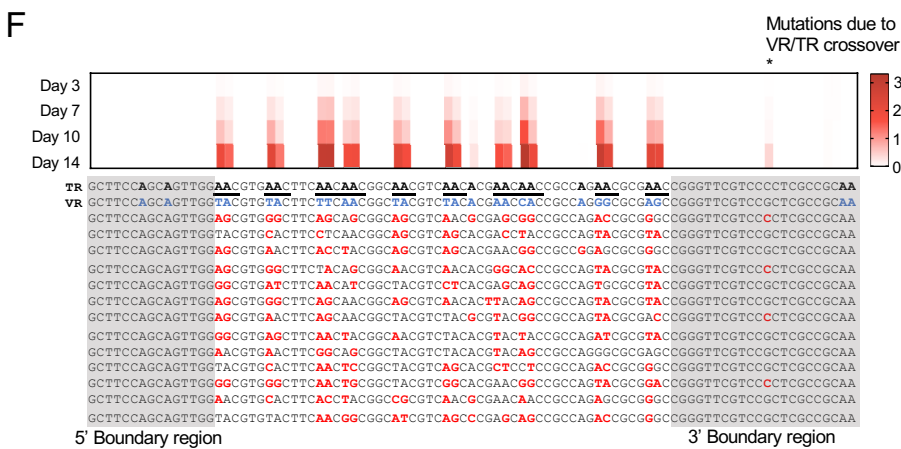
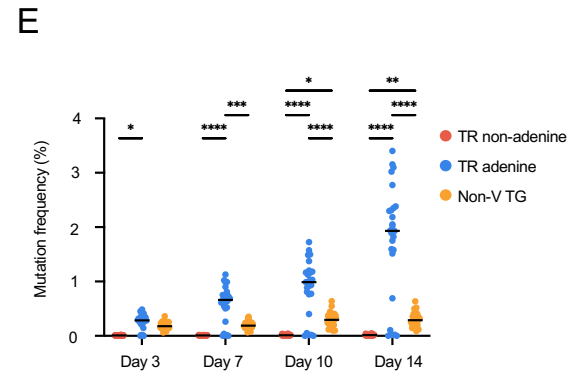
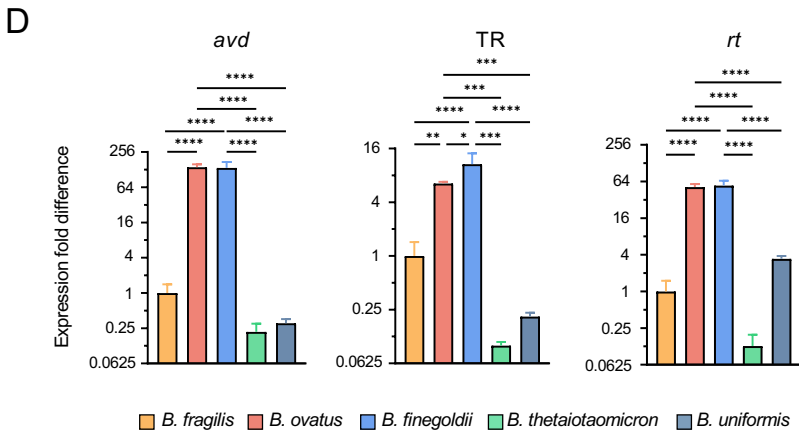
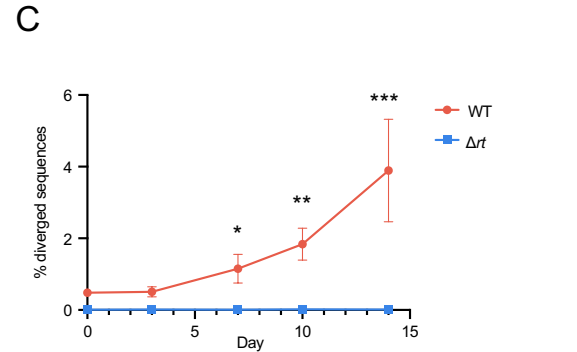
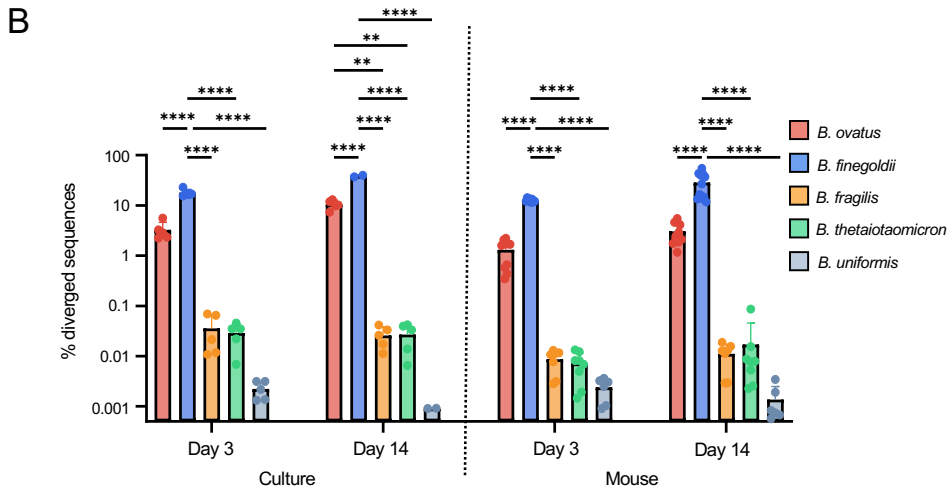
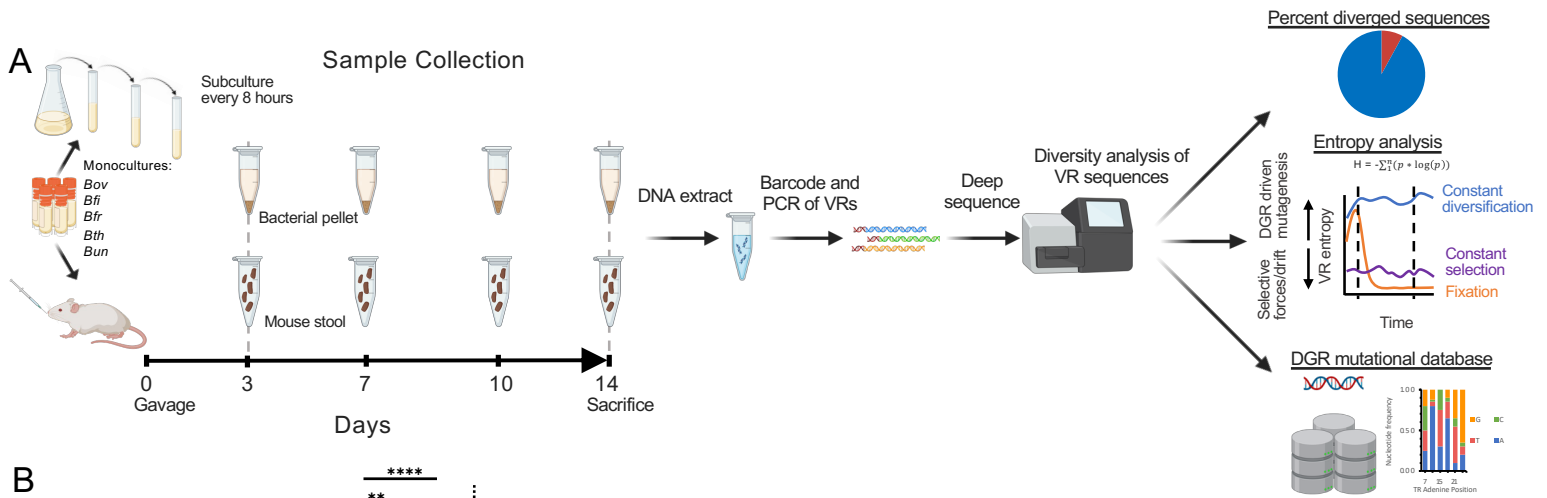
632 (E) The transfer efficiency, defined as the number of transconjugant cells divided by the total number
633 of recipient cells, of *in vitro* ICE mating assays.

634

635 Statistical analysis: *p<0.05, **p<0.01, ***p<0.001, ****p<0.0001. ANOVA test with Holm-Sidak
636 multiple comparison correction. Error bars, standard deviation of mean.

637

638 See also Supplemental Figure S4 and Supplemental Tables S9-13.



639 **Figure 4. *Bacteroides* DGRs are differentially active *in vitro* and *in vivo***

640 (A) Schematic of the experimental design.

641 (B) Percent of VRs that mutated from the parental VR sequence in cells grown *in vitro* or present in
642 stool samples from monocolonized SW mice (n=5-11). Error bars, standard deviation of mean.

643 (C) Percent of VR sequences that diverged from their cognate parental VRs in *Bov* WT and *Bov* Δrt
644 in stool samples from monocolonized SW mice (n=4). Error bars, standard deviation of mean.

645 (D) Expression of DGR encoded genes *avd*, TR, and *rt* measured by RNA-Seq and normalized to
646 *gyrA*, from mid-log phase cells grown *in vitro* (n=4). Error bars, standard deviation of mean fold
647 change.

648 (E) Mutation frequencies of individual nucleotide positions within the *Bov* target gene over a two week
649 period from a population of VRs extracted from stool samples of monocolonized mice (n=4).

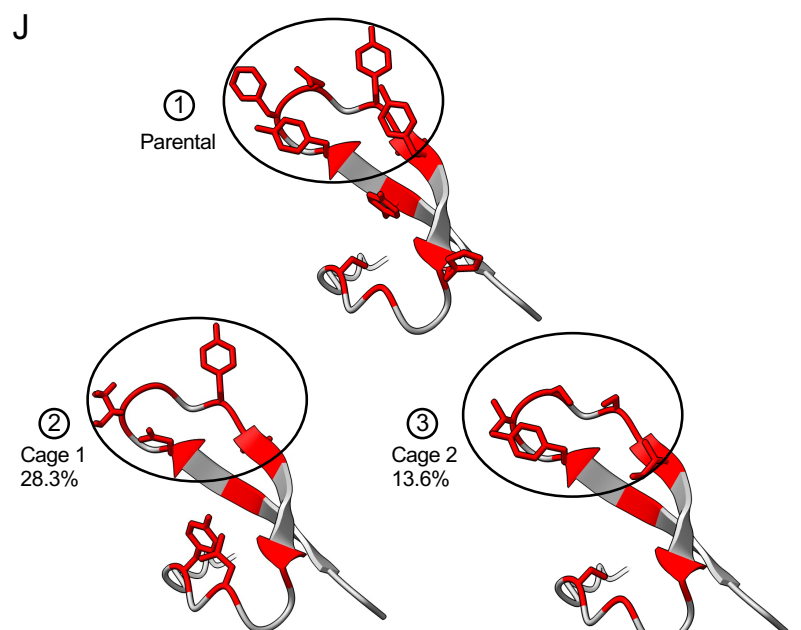
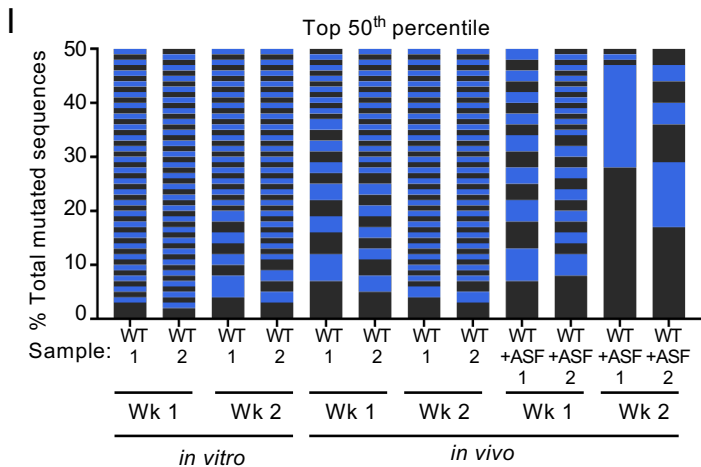
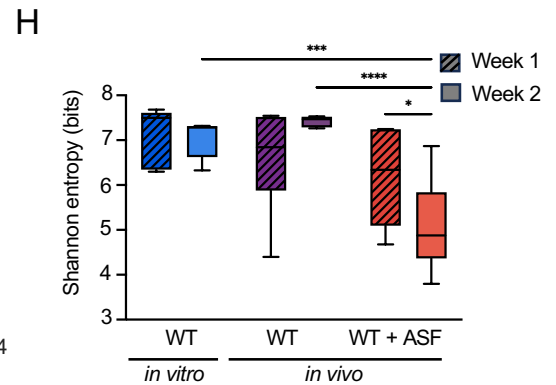
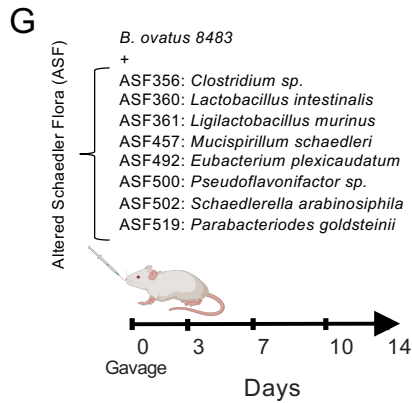
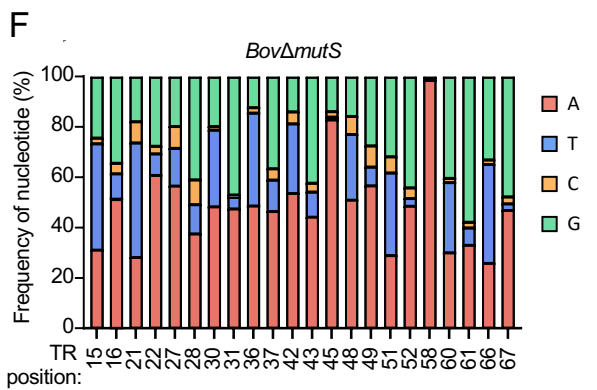
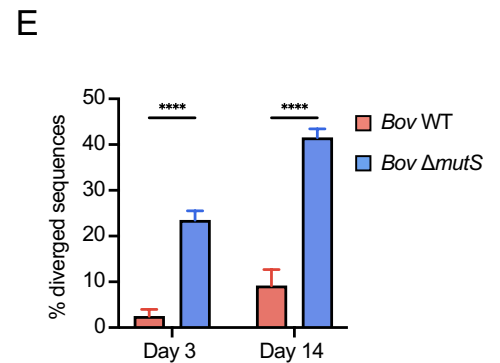
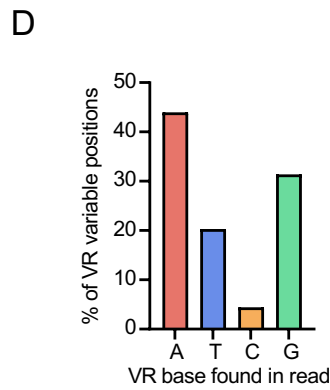
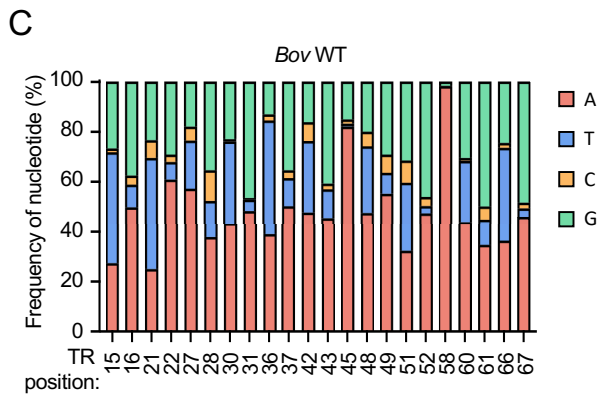
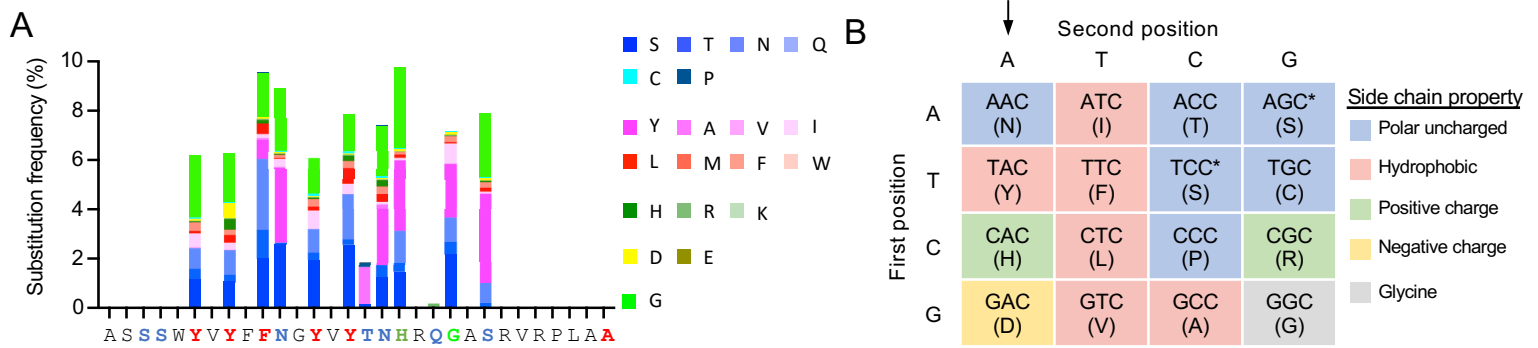
650 (F) Nucleotide entropy at each position along the *Bov* VR for a population of VRs extracted from stool
651 samples of monocolonized mice (n=4). Underneath the graph are the parental TR and VR
652 sequences followed by representative VR sequences at Day 14. TR adenines are bolded,
653 positions in VR corresponding to TR adenines are in blue, and positions in the mutated VR
654 sequence that are different from the parental VR are shown in red. The gray boxes show the
655 presumed 5' and 3' boundary regions where no adenine mutations are observed.

656 (G) Distribution of the number of nucleotide substitutions per VR read in *Bov* WT vs. *Bov* $\Delta mutS$ cells
657 grown *in vitro*.

658

659 Statistical analysis: *p<0.05, **p<0.01, ***p<0.001, ****p<0.0001. ANOVA test with Tukey multiple
660 comparison correction (B,D,G) or Holm-Sidak multiple comparison correction (C).

661 See also Supplemental Figures S5.



① A S S S W Y V Y F F N G Y V Y T N H R Q G A S R V R P L A A
 GCT · TCC · AGC · AGT · TGG · TAC · GTG · TAC · TTC · TTC · AAC · GGC · TAC · GTC · TAC · ACG · AAC · CAC · CGC · CAG · GGC · GCG · AGC · CGG · GTT · CGT · CCG · CTC · GCC · GCA

② A S S S W G V N F T G G Y V G T N G R Q N A Y R V R P L A A
 GCT · TCC · AGC · AGT · TGG · GGC · GTG · AAC · TTC · ACC · GGC · GGC · TAC · GTC · GGC · ACG · AAC · GGC · CGC · CAG · AAC · GCG · TAC · CGG · GTT · CGT · CCG · CTC · GCC · GCA

③ A S S S W G V Y F T S G S G I T N G R Q G A S R V R P L A A
 GCT · TCC · AGC · AGT · TGG · GGC · GTG · TAC · TTC · ACC · AGC · GGC · AGC · GGC · ATC · ACG · AAC · GGC · CGC · CAG · GGC · GCG · AGC · CGG · GTT · CGT · CCG · CTC · GCC · GCA

662 **Figure 5. *Bacteroides* DGRs are poised to alter amino acid side chain chemistry and respond to**
663 **competition *in vivo***

- 664 (A) Substitution frequency at *Bov* VR codons that correspond to TR adenines from cells grown *in*
665 *vitro*.
- 666 (B) Table showing codons that can be generated through mutagenic retrohoming of a TR AAC motif,
667 colored by the chemical class of the amino acid side chain.
- 668 (C) Nucleotide frequency at individual variable VR positions within diversified *Bov* VRs. The
669 corresponding TR adenine position is listed below each bar.
- 670 (D) Cumulative frequency of each of the four nucleotides at variable VR sites within diversified *Bov*
671 VRs.
- 672 (E) Percent of VRs that mutated from the parental VR sequence in *Bov* WT and *Bov* Δ *mutS* cells
673 grown *in vitro*.
- 674 (F) Nucleotide frequency at individual variable VR positions within diversified *Bov* Δ *mutS* VRs. The
675 corresponding TR adenine position is listed below each bar.
- 676 (G) Co-colonization of germ free mice with *Bov* with or without Altered Schaedler Flora (ASF), an 8
677 member bacterial community.
- 678 (H) Shannon entropy of diversified VR sequences obtained from *Bov* cells grown *in vitro* or present
679 in fecal samples of gnotobiotic mice (n=5-11). *p<0.05, **p<0.01, ***p<0.001, ANOVA with Tukey
680 multiple comparison correction.
- 681 (I) Graphical representation of diversified *Bov* VR populations derived from cells grown *in vitro* or
682 present within fecal samples of gnotobiotic mice. The height of the horizontal bars represents the
683 frequency of appearance, with alternating black and blue indicating unique VR sequences.
- 684 (J) VR encoded predicted structures of the parental *Bov* VR and two of the most commonly observed
685 VR sequences derived from fecal samples of separately caged gnotobiotic mice colonized with
686 *Bov* plus ASF bacteria. The VR nucleotide and amino acid sequences are shown below and the
687 frequency of the individual VR sequence within the diversified VR population is indicated.

688 Also see Supplemental Figure S6.

689 **Figure 6. DGRs undergo a burst of activity when transmitted from mother to infant**

690 (A) Number of unique DGRs found in mothers, in infants, or transmitted from mother to infant.

691 (B) Phylogeny of DGR RTs found in metagenomes derived from mothers, infants, and healthy

692 adults^{84,85,96,97}. Rings depict the genomic location of the DGR (ring 1), predicted host phyla (ring

693 2), variable protein homology (ring 3), and classification of the DGR host (ring 4). Red asterix

694 shows an area with Actinomycetota-harboring DGRs found primarily in infants.

695 (C) Number of DGRs identified per infant, grouped by age and mode of delivery. * $p < 0.05$, ** $p < 0.01$,

696 *** $p < 0.001$, ANOVA with Holm-Sidak multiple comparison correction.

697 (D) Number of DGRs identified per adult. n.s.: not significant.

698 (E) Taxonomic distribution of DGR-containing microbes, grouped by age and mode of delivery.

699 (F) Percent of DGRs with mutations at VR positions that correspond to TR adenines, grouped by age.

700 (G) Percent of DGRs where the predominant VR haplotype sequence changed between mother and

701 infant (n=388).

702 (H) Frequency of VR encoded amino acids from a DGR that was transmitted from mother to infant.

703 The target gene of this DGR is predicted to encode a Type V pilus tip adhesin.

704 (I) TR and VR sequences of *B. uniformis* 8492 (*Bun*) and the VRs from a similar DGR identified

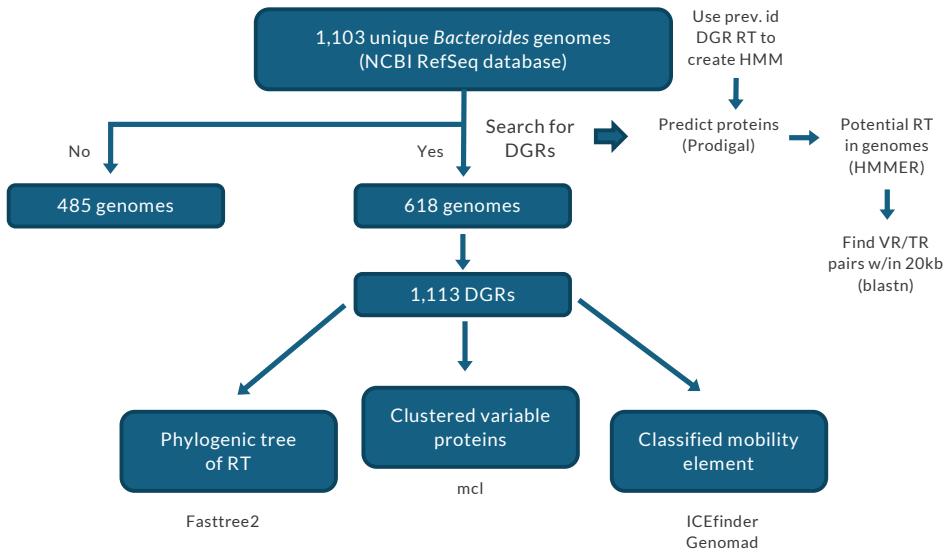
705 within a mother and her infant. Red, *Bun* TR adenines; Bolded, *Bun* VR nucleotides that

706 correspond to TR adenines; Blue, VR nucleotides that differ between *Bun* and the maternal VR

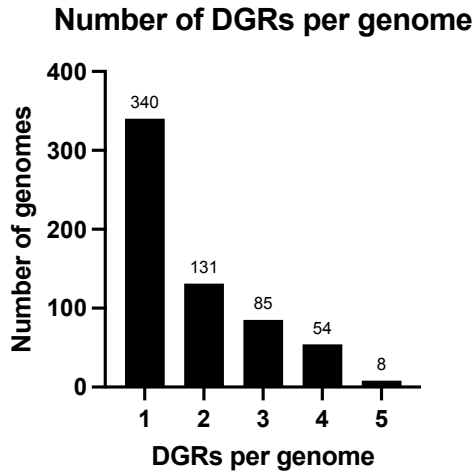
707 sequence; Orange, VR nucleotides that differ between the maternal and infant VR sequences.

708 See also Supplemental Figures S7 and Supplemental Tables S14-15.

A



B



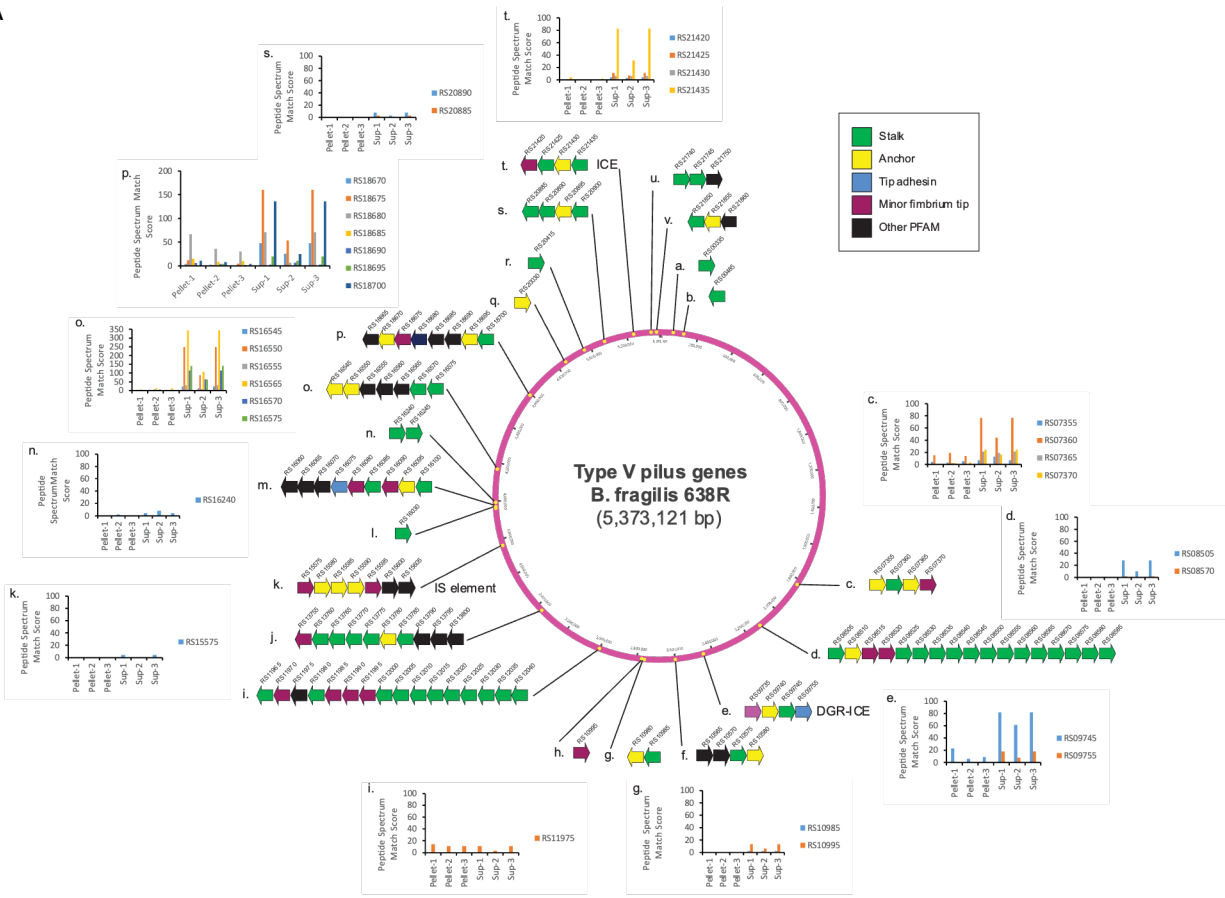
Supplemental Figure S1

709 **Supplemental Figure S1, Overview of *Bacteroides* DGRs, related to Figure 1:**

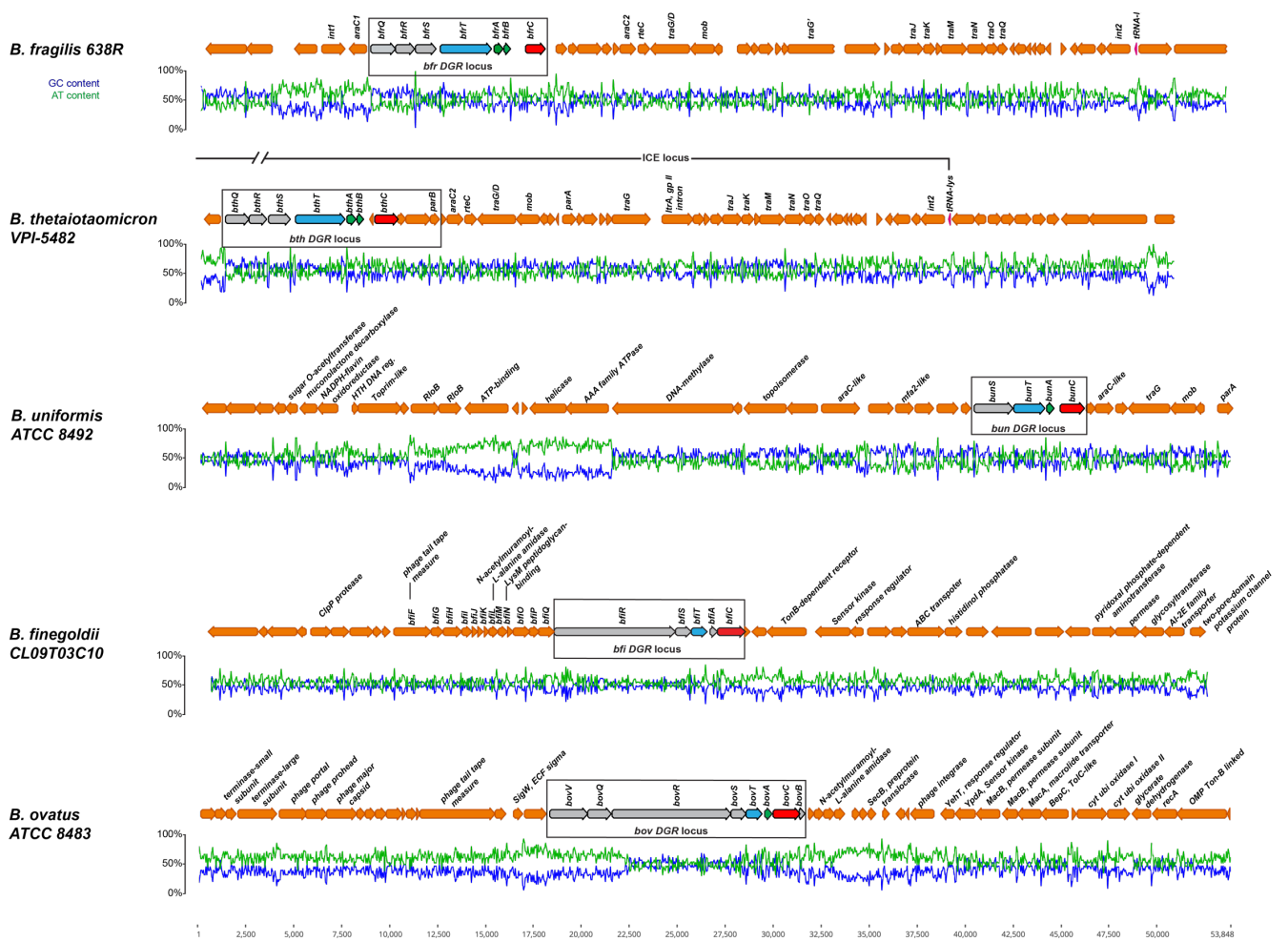
710 (A) Schematic of the DGR identification methodology from *Bacteroides* genomes. Proteomes were
711 first predicted from assemblies using Prodigal¹ and were used as input for profile-based
712 searches for the DGR RT proteins using HMMER². Imperfect repeats representing VR and TR
713 were searched within a 20 kb window upstream and downstream of potential *rt* genes using
714 blast³. Imperfect repeats were further filtered across two criteria. First, pairs had to differ from
715 each other at positions that correspond to adenines in one of the repeats. Second, the predicted
716 VR repeat had to be located within a gene encoding region previously predicted by Prodigal.
717 Using the identified DGRs, a phylogenetic tree was built using FastTree2⁴, the variable proteins
718 were clustered using blastp³ and mcl⁵. GeNomad⁶ and ICEBerg v2⁷ were used to determine if
719 the DGR locus fell within an ICE or prophage genome.

720 (B) Distribution of the number of DGRs identified within the same genome across all *Bacteroides*
721 strains in our dataset.

A



B



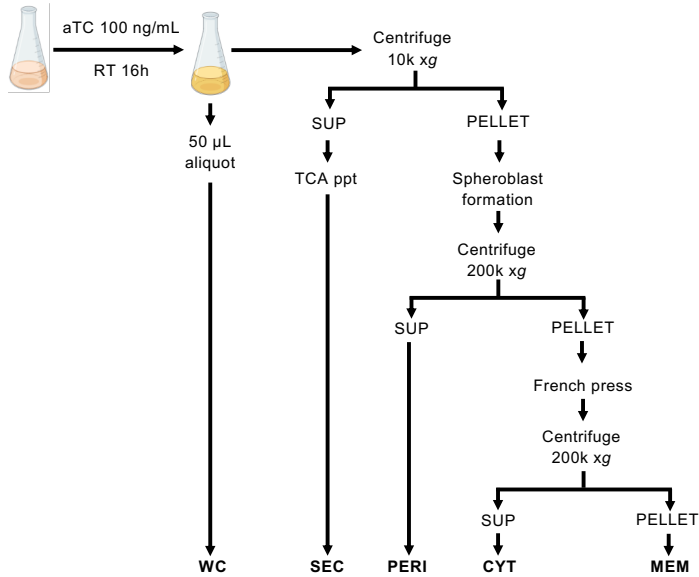
Supplemental Figure S2

722 **Supplemental Figure S2, *Bacteroides* pilin loci, related to Figure 2:**

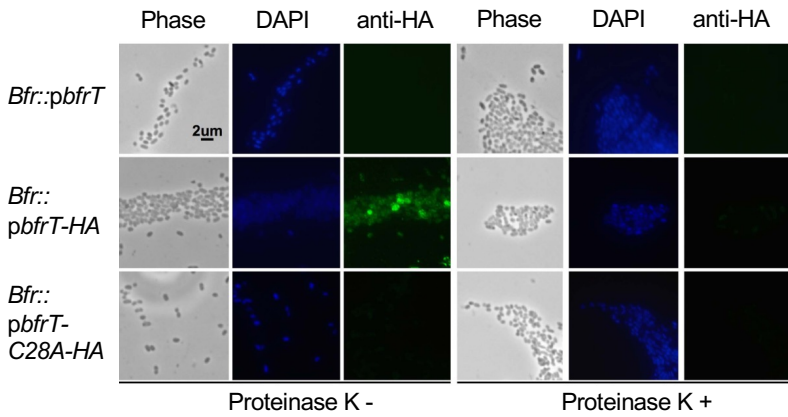
723 (A) Overview of the 22 loci within the *Bfr* genome that encode pilus proteins. Each gene was
724 functionally categorized according to its homology to known pilus proteins using HHpred^{8,9}. Bar
725 graphs for the indicated loci display the Peptide Spectrum Match Score across individual proteins
726 identified by mass spectroscopy of *Bfr* cells grown *in vitro*.

727 (B) Genomic loci and GC content adjacent to DGRs in the *Bacteroides* strains indicated. Within
728 each DGR locus, the target gene is colored blue, *rt* is red, and the accessory genes are green.
729 Predicted DGRs with accessory genes (colored gray) are boxed. Genes with predicted
730 annotations are indicated with the name of the gene. Genes without predicted names have
731 unknown function or are hypothetical.

A



C



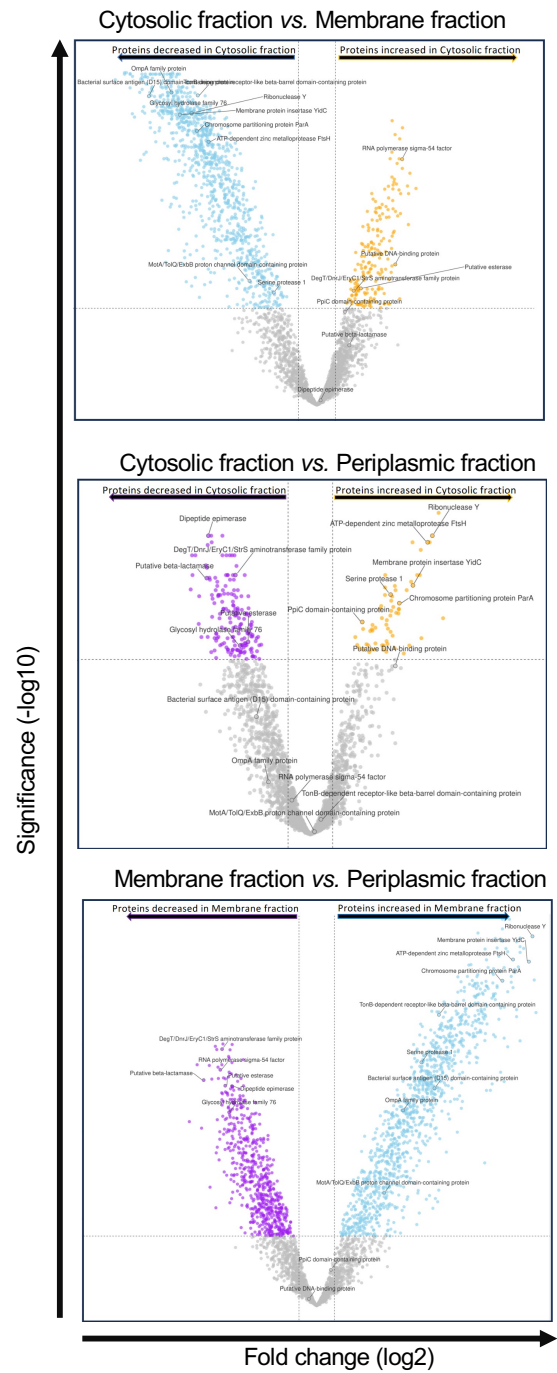
D

BfrT-ALFA Peptide Coverage

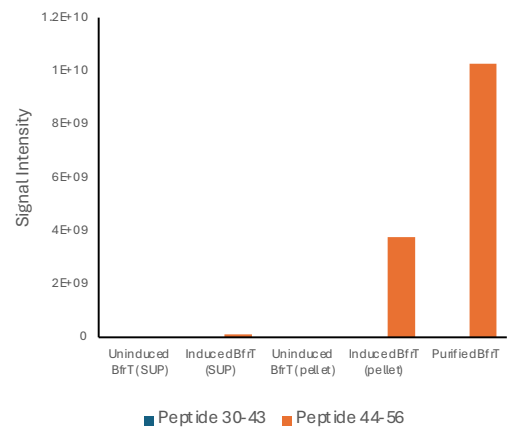
Red = Peptide mapped
Green = ALFA tag

10	20	30	40	50	60	70	80	90	100
MEMNKNLCRR	LGNLSLPVLL	SVVLLASCRD	EIETGAYTGP	YIRFSVSEGS	EWHSTRAAGG	PAEKAVPRDS	VQPLHGGDGN	TPLYLHTLYT	DSIASPSSDI
110	120	130	140	150	160	170	180	190	200
CPDTAVLTRA	TPVKTATLYE	SIGVLAAPFN	EPWSETSYRP	DYMYDVEVTK	ASSWTTSYHW	PTLTGGIRFF	AYAPYHGEI	VLSNKTGAGS	PTIITYVPAD
210	220	230	240	250	260	270	280	290	300
VADQKDLLFA	NSIYTTPTGT	LKNANNAAPL	TFNHALTAVR	FVCGNDMQEG	TVKSVSLNKG	CSKGLIINYGT	HSWSGVDTPA	DFSQTLDKST	TGTPDEALTT
310	320	330	340	350	360	370	380	390	400
DAQTFMMIPQ	TLPDGAQIEV	VFTDMSGTDY	TLTADIKGVV	WPIGKTVTYK	ISSSSINWTV	ELSVNMPGDF	TYSGGTQQYS	VTSYKHNSKG	DKQPAQWKAQ
410	420	430	440	450	460	470	480	490	500
FSEYGGPWID	TPPTWLTGFT	PFGAGGETSQ	SYNATVSAQI	GTSNDPHAQK	LRDNPShGGV	IYHHLNANQT	NGGSTDEMTA	NCYVVSQSGY	YCFPLYVYNA
510	520	530	540	550	560	570	580	590	600
IKNGAVNTSA	YTPFTGSGNI	LTTFFINHTGN	PITSPIYIKR	SGCVPAKAEI	LWQDAPGLIS	DVQYNSQMQ	LFVFNPNYIS	FQVNALTIKQ	GNAVIAIKDA
610	620	630	640	650	660	670	680	690	700
NDAILWSWHI	WVTDADINNV	IEVTNHQSQK	YKFMFVYLVG	CDGRTEYAE	RSCKVRFTAG	DASKEVLIKQ	VSASITVGGN	HPYEWGRKD	PFPPSDGLSN
710	720	730	740	750	760	770	780	790	800
TNKTWYDKG	NAHTESPKTE	NFSTGATCIM	NYILKPDVMQ	SQFYGDNTYA	NLWSADNNVY	TANDENVIKT	IYDFSPVQFK	LPVGNFTTFG	TTTGNNTSTF
810	820	830	840	850	860	870	880	890	900
SEINGTWSS	LKGWNFPTDA	SRSKTIFFPA	SGYRVCSTGG	AANVGSYGSC	WSAVPHNQY	GRNLAFNSSN	VYPLNSDRA	YGFGRSSQE	SRLEEELRLR
EELRLRLTE									

B



E



732 **Supplemental Figure S3, *Bacteroides* diversify pilus proteins, related to Figure 2**

733 (A) Schematic of cellular fractionation methodology. Cells were either induced with 100 ng/mL of aTC
734 or not and incubated at room temperature for 16 hours. A 50 μ L aliquot was set aside for the
735 whole cell fraction. The remaining culture was centrifuged at 10,000 x g. The supernatant was
736 TCA precipitated and stored as the secreted fraction. The pellet fraction was subjected to
737 spheroblasting and centrifuged at 200,000 x g. The resulting supernatant was stored as the
738 periplasmic fraction while the pelleted fraction was passed through a French press to separate
739 cytoplasmic and crude membrane fractions.

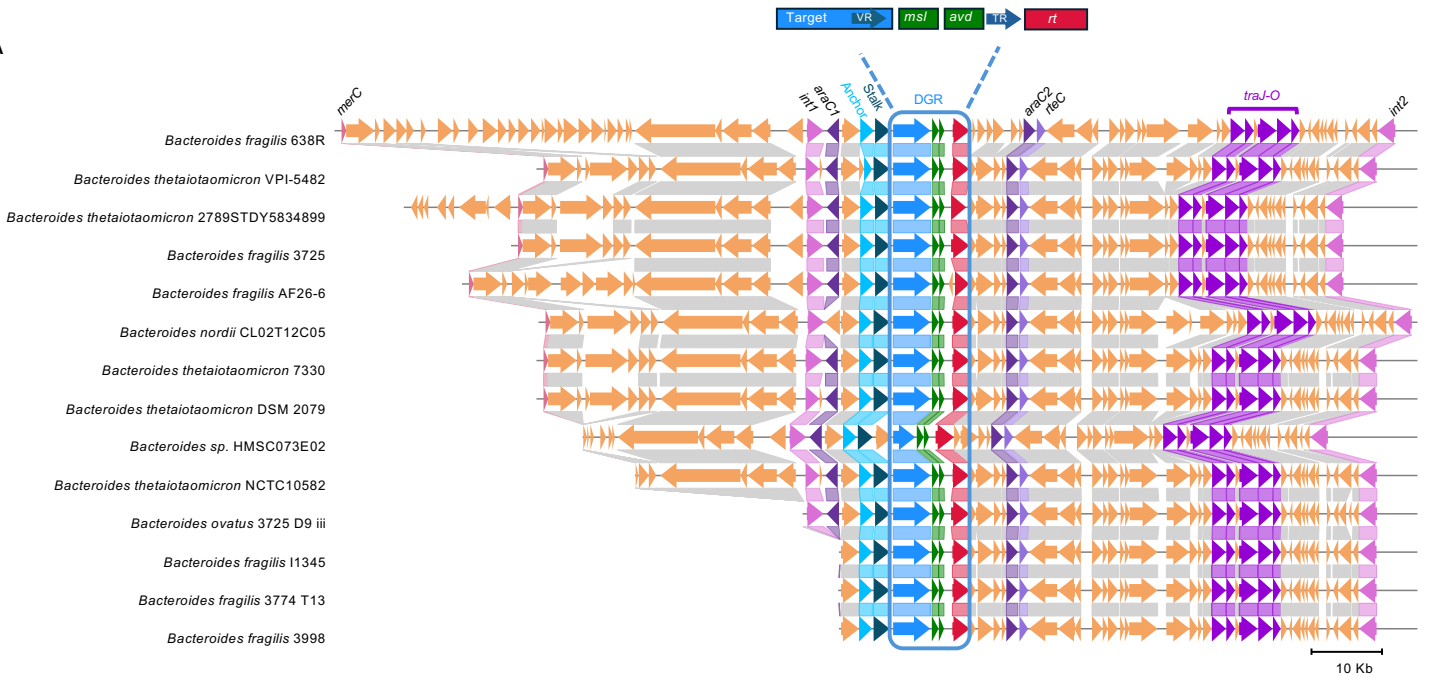
740 (B) Volcano plot of relative protein enrichment between cellular fractions in *Bov* cells grown *in vitro*.
741 The abundance of individual proteins from the periplasmic, cytosolic, and membrane fractions
742 was quantified by mass spectrometry and the relative abundance change of each protein was
743 calculated for each fraction comparison (i.e., cytosol fraction compared to membrane fraction;
744 cytosol fraction compared to periplasmic fraction; membrane fraction compared to periplasmic
745 fraction). The VolcanoR¹⁰ web application was used to plot the Log₂ fold change of protein
746 abundance versus the -Log₁₀ of the Significance value of the change in abundance for each
747 protein between each of the fractions. Specific proteins with previously reported sub-cellular
748 localization data were tracked to evaluate the effectiveness of the fractionation protocol (see
749 Supplemental Table S8).

750 (C) Immunofluorescence of *Bfr* cells expressing epitope tagged BfrT. Cells overexpressing WT BfrT
751 (untagged), HA-tagged BfrT, or the HA-tagged BfrT-C28A mutant were treated with or without
752 proteinase K.

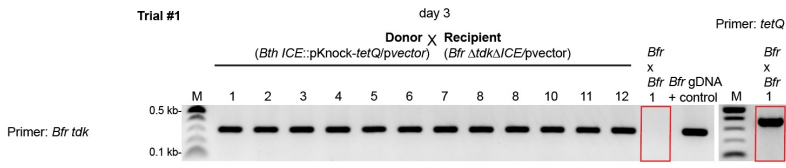
753 (D) Summary of the peptides identified by mass spectroscopy that aligned to BfrT from the
754 supernatant of *Bfr* cells grown *in vitro*.

755 (E) Signal intensity of the peptide fragments recovered from Part D for peptides 30-43 and 44-56,
756 which represent arginine cleavage sites during pilus assembly.

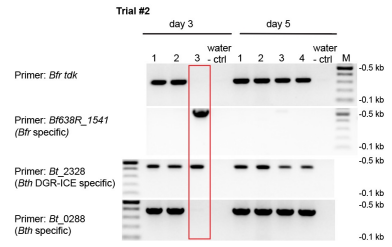
A



B



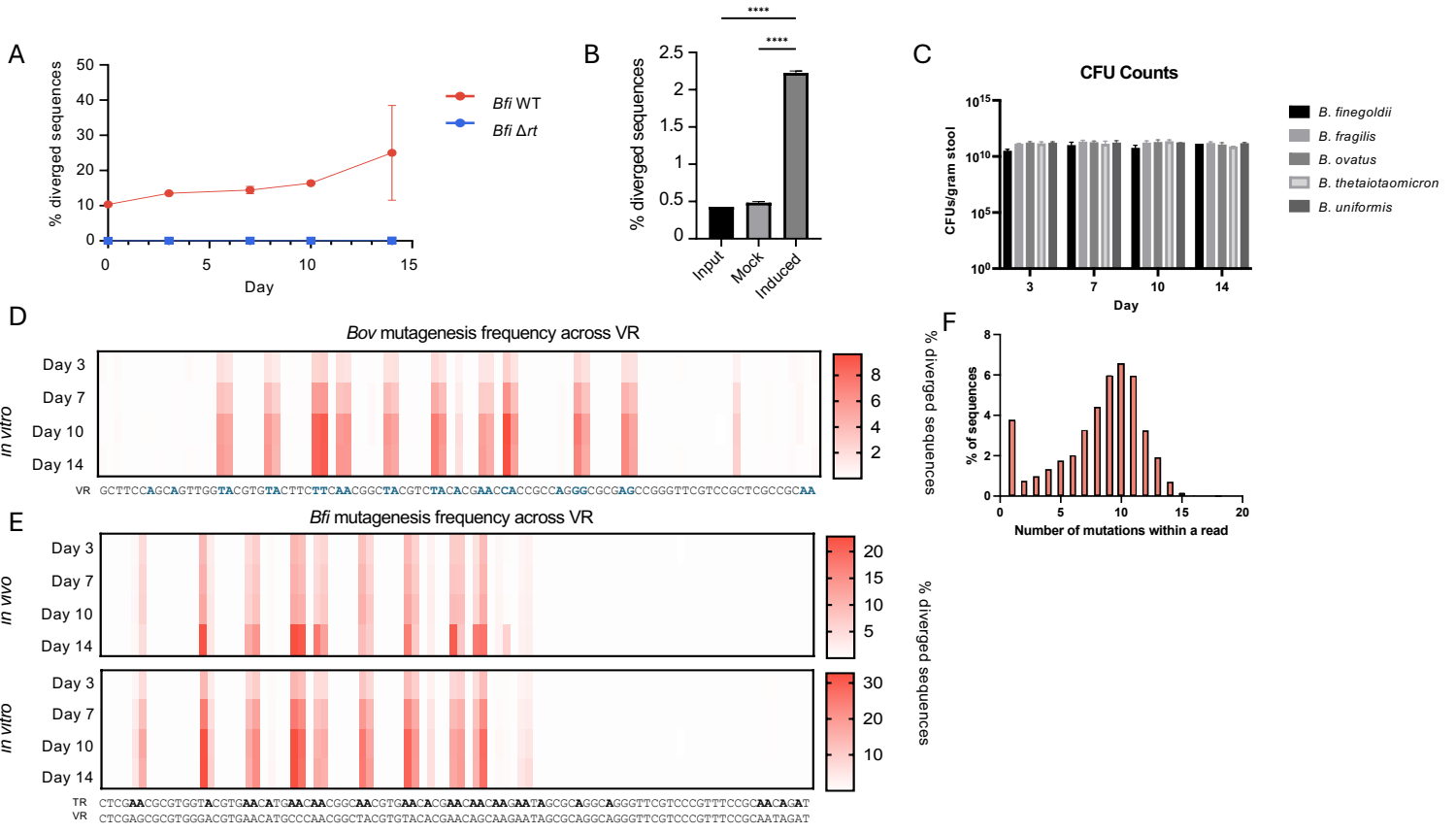
C



Supplemental Figure S4

757 **Supplemental Figure S4, Overview of *Bacteroides* ICEs, related to Figure 3**

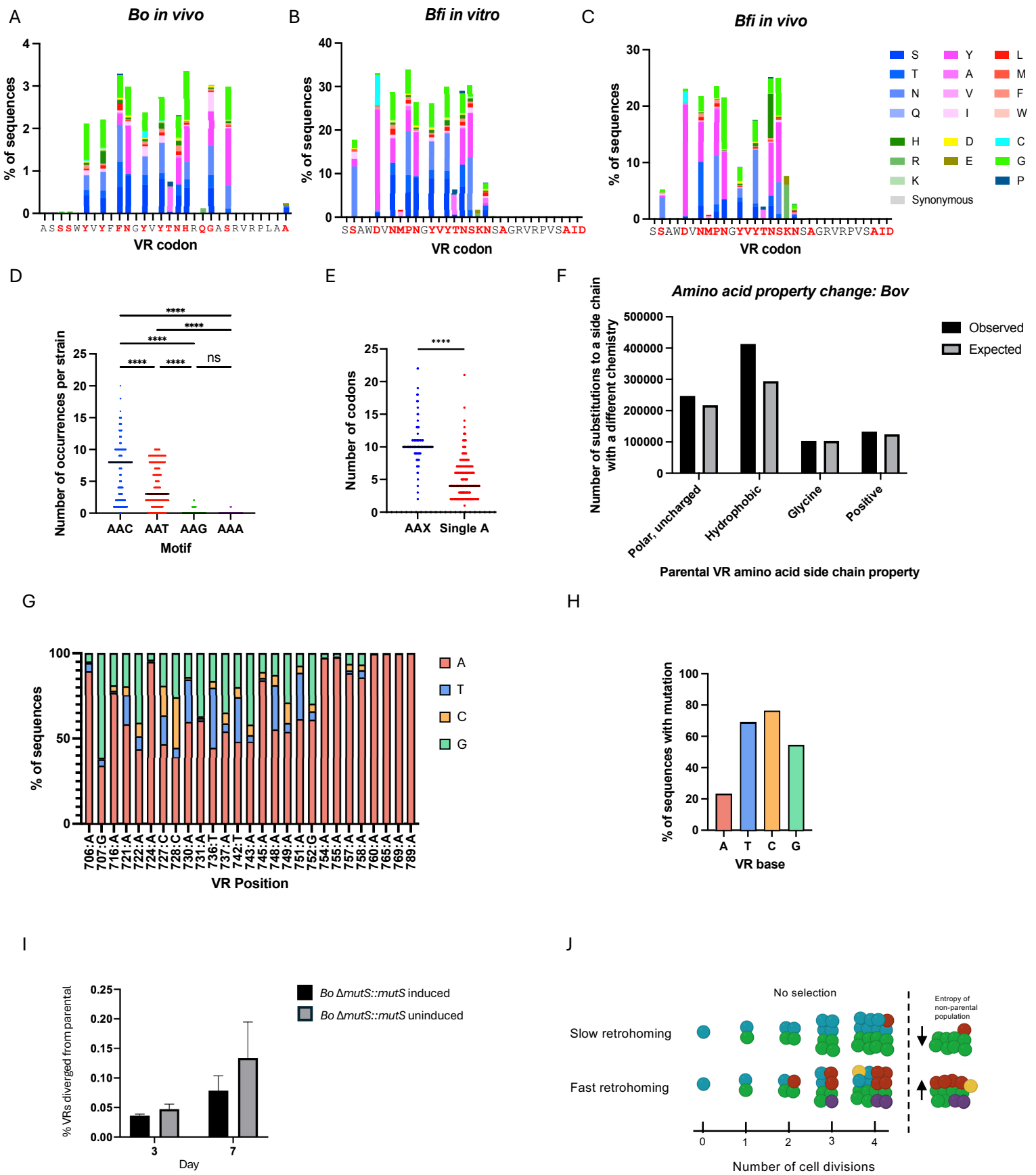
- 758 (A) Synteny of 16 DGR-encoding ICEs identified in different *Bacteroides* strains. The name of the
759 *Bfr* gene is labeled on top for genes with predicted functions. Genes are connected if their
760 encoded proteins have at least 90% identity. The DGR locus within the ICEs is boxed.
- 761 (B) PCR based screen to identify ICE transconjugants from matings between *Bfr* and *Bth* cells. True
762 transconjugants, as shown by the red box, are negative for the *tdk* gene while positive for the
763 ICE (*tetQ*).
- 764 (C) PCR based screen as in Part B but using *Bth* cells as donor cells.



Supplemental Figure S5

765 **Supplemental Figure S5, DGRs are active in *Bacteroides*, related to Figure 4**

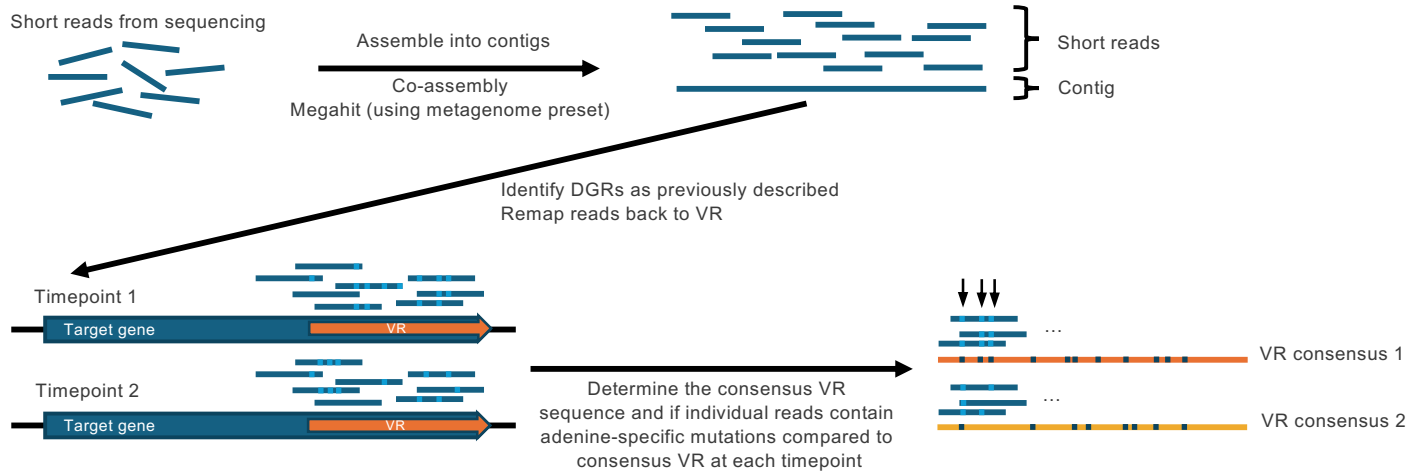
- 766 (A) Percent of VRs that diverged from their cognate parental VR in *Bfi* WT and *Bfi* Δrt present in
767 fecal samples from monocolonized SW mice (n=4 each). Error bars, standard deviation of
768 mean.
- 769 (B) Number of VR reads that diverged from their cognate parental VR sequence in *Bov* $\Delta rt::pNBU2-$
770 *tetR* P1T_{DP}^{GH023}-*rt*, a mutant in which the chromosomal copy of the DGR *rt* has been knocked out
771 and WT *rt* is expressed ectopically under the control of an inducible promoter. Mock induction or
772 induction of *rt* transcription with aTC are shown. Error bars, standard deviation from mean.
773 ****p<0.0001, ANOVA.
- 774 (C) Number of CFUs obtained per gram of stool from monocolonized SW mice over a two-week
775 period (n=4); error bars, standard deviation from mean. All comparisons across strains are non-
776 significant (n.s.).
- 777 (D) Mutation frequencies of individual nucleotide positions within the *Bov* target gene over a two-
778 week period from cells grown *in vitro*.
- 779 (E) Mutation frequencies of individual nucleotide positions within the *Bfi* target gene over a two-week
780 period present in fecal samples from monocolonized SW mice (top) or from cells grown *in vitro*
781 (bottom).
- 782 (F) Distribution of the number of nucleotide substitutions per VR read in *Bfi* cells grown *in vitro* at
783 Day 14.



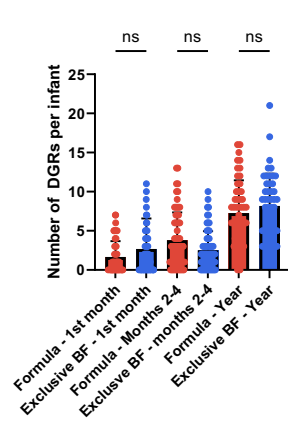
784 **Supplemental Figure S6, *Bacteroides* DGRs preferentially create non-synonymous mutations,**
785 **related to Figure 5**

- 786 (A) Substitution frequency at *Bov* VR codons present in fecal pellets of monocolonized SW mice at
787 14 days post-gavage.
- 788 (B) Substitution frequency at *Bfi* VR codons from cells grown *in vitro* after 14 days.
- 789 (C) Substitution frequency at *Bfi* VR codons present in fecal samples of monocolonized SW mice at
790 14 days post-gavage.
- 791 (D) Number of AAX motifs present in *Bacteroides* TRs. **** $p < 0.0001$, ANOVA, Holm-Sidak multiple
792 comparison correction. N.s, not significant by ANOVA test.
- 793 (E) Number of AAX motifs vs motifs containing a single adenine in *Bacteroides* TRs. **** $p < 0.0001$,
794 Mann-Whitney test.
- 795 (F) Number of occurrences that a diversified codon would encode for an amino acid with a different
796 chemical property of its side chain from *Bov* cells grown *in vitro*. The expected number was
797 calculated under the assumption that every AAX motif has an equal chance of occurring through
798 mutagenic retrohoming.
- 799 (G) Nucleotide frequency at individual variable VR positions within diversified *Bfi* VRs. The parental
800 VR position and nucleotide are listed under each bar.
- 801 (H) Cumulative frequency of each of the four nucleotides at variable VR sites within diversified *Bfi*
802 VRs.
- 803 (I) Percent of VRs that diverged from their cognate parental VR in *Bov* Δ *mutS*::*pNBU2-tetR*
804 P1T_{DP}^{GH023}-*mutS*, a mutant in which the chromosomal copy of *mutS* was knocked out and WT
805 *mutS* was expressed ectopically under an aTC-inducible promoter, following aTC induction or
806 mock induction (n=3). Error bars, standard deviation from mean.
- 807 (J) Representation of the effect on VR diversity of slow mutagenic retrohoming vs. fast mutagenic
808 retrohoming with respect to cell division. Blue spheres represent cells with the parental VR
809 sequence, all other colored spheres represent cells with diversified VR sequences.

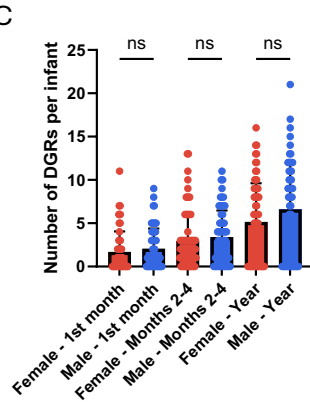
A



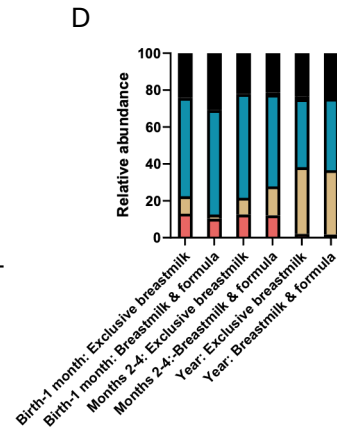
B



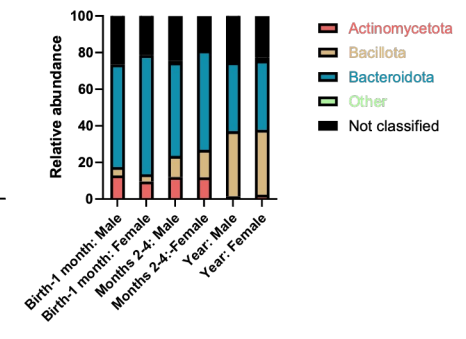
C



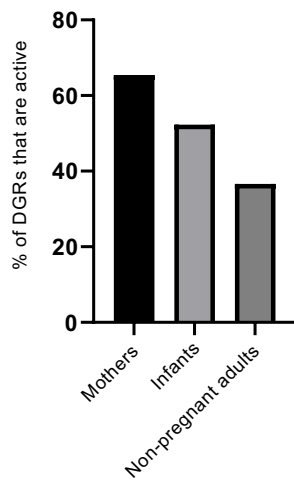
D



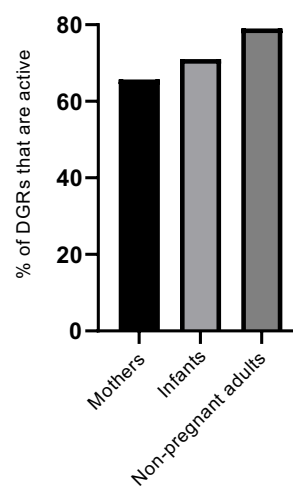
E



F

DGRs similar to *Bfr*, *Bth*, and *Bun*

G

DGRs similar to *Bov* and *Bfi*

810 **Supplemental Figure S7, DGRs in mother-infant datasets, related to Figure 6**

- 811 (A) Schematic of the methodology for DGR identification and activity detection from metagenomic
812 sequencing reads.
- 813 (B) Number of DGRs identified per infant, grouped by feeding preference and age. n.s. not significant
814 by ANOVA tests.
- 815 (C) Number of DGRs identified per infant, grouped by gender and age. n.s. not significant.
- 816 (D) Taxonomic distribution of DGR-containing microbes at different ages, grouped by feeding
817 preference.
- 818 (E) Taxonomic distribution of DGR-containing microbes at different ages, grouped by gender.
- 819 (F) Percent of active DGRs with target genes similar to the target genes of *Bfr*, *Bth*, and *Bun* in
820 mothers, infants, and healthy adults.
- 821 (G) Percent of active DGRs with target genes similar to the target genes of *Bov* and *Bfi* in mothers,
822 infants, and healthy adults.

823 **Methods**

824

825 **Bacterial strains, plasmids, primers, and growth conditions.** Bacterial strains are described in
826 Supplementary Table 1, plasmids are described in Supplementary Table 2, and primers are described in
827 Supplementary Table 3. *E. coli* strains were grown aerobically at 37°C in Luria-Bertani medium (LB, BD
828 Difco). *Bacteroides* strains were grown in Brain Heart Infusion (BHI, BD Difco) medium at 37°C,
829 supplemented with vitamin K (5 µg/mL) and hemin (5 µg/mL), and incubated in an anaerobic chamber
830 (Coy) under 5% H₂, 10% CO₂, 85% N₂. Antibiotics were added to the following final concentrations:
831 Carbenicillin (Cb), 100 µg/mL; Gentamicin (Gm), 200 µg/mL; Erythromycin (Erm), 5 or 25 µg/mL,
832 Tetracycline (Tet), 2 µg/mL, 5-fluoro-2'-deoxyuridine (FudR, 200 µg/mL), unless otherwise noted.
833 Anhydrous tetracycline (aTC) was used at a concentration of 100 ng/mL. Backbone plasmids used in this
834 study were pLGB13, pLGB36, and pNBU2_erm-TetR-P1T_DP-GH023^{102–104}.

835 **Mice.** Germ-free Swiss Webster mice were purchased from Taconic Farms and bred in flexible film
836 isolators. For gnotobiotic experiments, sterile litters of 8-10 week old male and female mice were
837 transferred into autoclaved microisolator cages where they were fed autoclaved chow diets *ad libitum*
838 and given autoclaved water supplemented with gentamicin 100 µg/mL. Altered Schaedler Flora (ASF)
839 live animal donor C57BL/6 mice were purchased from Taconic (ASF-DONOR-M/F) and fed an autoclaved
840 chow diet *ad libitum* and given autoclaved water. For monocolonization experiments, animals were orally
841 gavaged with a 200 µL inoculum of an overnight culture of a given *Bacteroides* strain containing ~10⁸
842 cells per gavage. For ASF experiments, approximately 1 g of fresh stool from an ASF donor mouse was
843 vortexed into 1 mL of an overnight culture of *Bacteroides* immediately prior to oral gavage. For
844 conjugation experiments, 100 µL of each donor and recipient strains were combined in a 1:1 ratio and
845 orally gavaged to mice. Stool samples were collected on Days 3, 7, 10, and 14. All procedures were
846 performed in accordance with an approved protocol following IUCAC guidelines at the University of
847 California, Los Angeles and the IUCAC guidelines at the California Institute of Technology.

848 **Plasmid generation.** Q5 high fidelity DNA polymerase (New England Biolabs [NEB]) or Phusion high
849 fidelity DNA polymerase (NEB) were used for PCR cloning steps. To construct plasmids for deletion
850 mutagenesis, a PCR fragment was generated that included 1kb upstream and the first 12 codons of the
851 targeted gene and a second fragment was generated to include the last 12 codons and 1kb downstream
852 of the targeted gene. The backbone plasmids pLGB36 (*Bfr* only) or pLGB13 (all others) were cut at the
853 BamHI (NEB) restriction site. PCR fragments were ligated together and into the backbone plasmid using
854 NEBuilder HiFi assembly master mix (NEB). For overexpression mutants, the targeted gene was cloned
855 into pNBU2_erm-TetR-P1T_DP-GH023 plasmid cut at NcoI and Sall restriction sites using the NEBuilder

856 HiFi assembly master mix. Plasmids were sequenced following cloning to ensure insertions were without
857 mutations.

858 **Conjugation and allelic exchange.** Plasmids were transformed into *E. coli* S17- λ pir and conjugated into
859 *Bacteroides* strains as previously described¹⁰². Briefly, *E. coli* S17- λ pir donor cultures were grown in 25
860 mL of LB+Cb to an OD₆₀₀ between 0.3-0.5. Recipient *Bacteroides* strains were grown in BHIS to an OD₆₀₀
861 of 0.05-0.1. Both strains were mixed together in a 50 mL Falcon centrifuge tube and spun at 4000xg for
862 10 mins. Cell pellets were then resuspended in 100 μ L of BHIS, plated on the center of a prewarmed
863 BHIS plate, and incubated aerobically at 37°C for 14-16 hours. The mating spot was resuspended in
864 BHIS and was diluted serially from 1:10 to 1:10,000. From each dilution, 100 μ L was streaked onto
865 BHIS+Gm+Erm selection plates and incubated at 37°C anaerobically. Colonies were picked after two
866 days, restreaked on BHIS+Gm+Erm plates, and grown for two additional days. For overexpression
867 mutants, stocks were made from an overnight culture of BHIS and integration was confirmed by PCR of
868 attB sites using the appropriate primers in Supplemental Table 3. For allelic exchange protocols, isolates
869 were grown in BHIS overnight, diluted to 1:1,000 or 1:10,000 in the morning, and plated on BHIS+aTC100
870 to induce the counterselection toxin. After two days, colonies were picked, and colony PCR screening
871 was performed to determine which colonies contained the desired mutation and which colonies reverted
872 back to WT. Stocks were made of potential mutants. The mutation site including 1kb upstream and
873 downstream of the mutation, as well as the entire DGR locus was amplified from chromosomal DNA and
874 sequenced to ensure no undesired mutations were introduced.

875 **Sample preparation and immunoblotting.** For SDS-polyacrylamide gel electrophoresis (SDS-PAGE)
876 sample preparation, *Bacteroides* strains were cultured in BHIS media and harvested as described
877 previously¹⁰⁵. Proteins with ALFA-tags were detected using a monoclonal mouse antibody at a dilution of
878 1:5000 (anti-ALFA, NanoTag Biotechnologies). Immunodetection was carried out by chemifluorescence
879 using horseradish peroxidase-labelled goat anti-mouse IgG and the ECL plus[®] detection substrate (GE
880 Healthcare). Chemifluorescent signals were visualized using a Typhoon scanner (GE Healthcare).

881 **Cellular fractionation.** Cellular fractionation was carried out as described¹⁰⁶ and is briefly summarized
882 here. Overnight cultures of WT or mutant *Bacteroides* strains encoding inducible ALFA-tagged proteins
883 were diluted 1:100 in a total volume of 1 L BHIS supplemented with 25 μ g/ml Erm and incubated at 37°C
884 until reaching OD 0.05-0.1. The target gene was then induced by addition of 100 ng/mL of aTC to the
885 media and the culture was allowed to incubate for an additional 16 hr at RT. The induced 1 L culture was
886 then pelleted by centrifugation at 10,000xg and subsequently resuspended in 10 mL of spheroblast buffer
887 (0.2 M Tris-HCl pH 8, 1 M sucrose, 1 mM EDTA, 1 mg/mL lysozyme) and incubated for 5 min at RT. A

888 volume of 40 mL of ice-cold dH₂O was then added to the suspension before placing on ice for 5 min to
889 allow spheroblast formation. The suspension was then centrifuged at 200,000xg for 45 min at 4°C. The
890 resulting supernatant was collected as the periplasmic fraction and the pellet was resuspended in French
891 press buffer (7.5 mL ice-cold 10 mM Tris-HCl pH 7.5, 5 mM EDTA, 0.2 mM DTT, 50 µL 1mg/mL DNaseI).
892 Cells were then ruptured in a French Press with two passes at 10⁸ Pa. Unbroken cells in the lysate were
893 removed by centrifugation at 10,000xg for 10 min at 4°C. The lysate was then centrifuged at 280,000xg
894 for 4hr at 4°C. The resulting supernatant was collected as the cytoplasmic fraction while the pellet
895 contained crude membranes. Membrane fractions were diluted 1:1 with dH₂O, centrifuged at ≥85 000xg
896 for 20 min at 4°C, washed 3x in 500 µL dH₂O and then stored at -20°C.

897 **Proteinase K.** Overnight cultures of *Bacteroides* strains carrying an inducible, ALFA-tagged variable
898 protein (*Bfr* WT-ALFA, *Bfr* C28A-ALFA, or *Bov* WT-ALFA) were diluted 1:100 in a total volume of 10 mL
899 BHIS supplemented with 25 µg/ml Erm and incubated at 37°C until reaching OD 0.1. Expression of the
900 ALFA-tagged variable protein was then induced by addition of 10 ng/mL aTC. Induced cultures were
901 incubated for 8 hrs and then collected by centrifugation. Cells were resuspended at 2.5 OD/mL in a total
902 volume of 5 mL PBS and 1 mL aliquots of the bacterial slurry were dispensed to 1.5 ml tubes and
903 incubated for 1 hour at 37°C with one of the following Proteinase K quantities : 0 (control); 25 ng; 50 ng;
904 100 ng; 200 ng. Next, 10 µl PMSF was added to each tube and cells were collected by centrifugation at
905 10,000 xg and washed 2X with 1 mL PBS. Cells were resuspended in SDS sample dye with beta-
906 mercaptoethanol, boiled for 10 minutes and analyzed by Western blot.

907 **ALFA pulldown.** ALFA-tagged *Bacteroides* variable proteins and their interacting protein partners were
908 purified from an appropriate cellular fraction (periplasm for *Bov*, supernatant for *Bfr*) using the Anti-ALFA
909 single domain nanobody resin (ALFA SelectorST) according to manufacturer instructions. Briefly, 3 mL
910 of the cellular fraction was diluted 1:1 with 2x binding buffer (50 mM Tris pH 7.5, 100 mM NaCl, 2 mM
911 EDTA, 1 % NP-40, 10% glycerol) supplemented with 10 uL/ml HALT protease (ThermoFisher). To this
912 suspension, 200 µL of ALFA SelectorST was added and the mixture was incubated at 4°C with end-over-
913 end rotation for 16 hr. The resin was collected by centrifugation and washed 2x with PBS containing 0.5%
914 NP-40. The resin was then pelleted and resuspended in 200 µL of 1x Laemmli buffer (0.0625 M Tris pH
915 6.8, 2% sodium dodecyl sulfate, 10% glycerol, bromphenol blue) with 1% β-mercaptoethanol and then
916 boiled for 5 min.

917 **Mass spectroscopy.** Immunoprecipitation eluates (ALFA pulldowns) or cellular fractions (fractionation
918 controls) in 1x Laemmli buffer were diluted in equal volume of 100mM Tris-Cl pH 8.5 and reduced and
919 alkylated by the sequential addition of 5 mM tris(2-carboxyethyl) phosphine and 10 mM iodoacetamide.

920 This was followed by treatment with single-pot, solid-phase-enhanced sample preparation (SP3) protocol
921 for protein clean-up¹⁰⁷. Following SP3, eluates were proteolytically digested with Lys-C and trypsin at
922 37°C overnight. The digested peptides were subjected to offline SP3-based peptide clean-up and
923 subsequently analyzed by LC-MS/MS. Briefly, peptides were separated by reversed-phase
924 chromatography using 75 µm inner diameter fritted fused silica capillary column packed in-house to a
925 length of 25 cm with bulk 1.9 mM ReproSil-Pur beads with 120 Å pores. The increasing gradient of
926 acetonitrile was delivered by a Dionex Ultimate 3000 (Thermo Scientific) at a flow rate of 200 nL/min.
927 MS/MS spectra were collected using data-dependent acquisition on an Orbitrap Fusion Lumos Tribrid
928 mass spectrometer (Thermo Fisher Scientific) with an MS1 resolution (r) of 120,000 followed by
929 sequential MS2 scans at a resolution (r) of 15,000. The data generated by LC-MS/MS were analyzed
930 using the MaxQuant bioinformatic pipeline¹⁰⁸. The Andromeda integrated in MaxQuant was employed as
931 the peptide search engine. Briefly, a maximum of two missed cleavages was allowed. The maximum
932 false discovery rate for peptide and protein was specified as 0.01. Label-free quantification (LFQ) was
933 enabled with LFQ minimum ratio count of 1. The parent and peptide ion search tolerances were set as
934 20 and 4.5 ppm respectively. The MaxQuant output files were subsequently processed for statistical
935 analysis of differentially enriched proteins using Analytical R tools for mass spectrometry (artMS)¹⁰⁹.

936 **Cellular Fractionation Controls.** To assess the quality of protein enrichment from our cell fractionation
937 protocol, we used mass spectrometry to identify all proteins in each fraction generated and then
938 compared the change in the abundance of each identified protein between each fraction (i.e., cytosol
939 fraction compared to membrane fraction; cytosol fraction compared to periplasmic fraction; membrane
940 fraction compared to periplasmic fraction). For each protein, the change in its abundance between
941 fractions was graphed using the VolcanoR¹¹⁰ web application. The log₂ fold change of abundance was
942 graphed over the x-axis and the -log₁₀ significance value of the change in abundance was graphed over
943 the y-axis. We then tracked specific proteins with previously reported sub-cellular localization data to
944 evaluate the effectiveness of the fractionation protocol (Supplementary Table S8).

945 **Detection of DGR-ICE integration, excision, and episome formation by PCR.** Genomic DNA was
946 prepared from cultured *Bfr* or *Bth* cells carrying overexpression constructs (pEV, *paraC2*, *prteC*, *pmerR*)
947 using a DNeasy Blood and Tissue Kit (Qiagen). A 50 ng aliquot of DNA was used as template in a
948 standard 50 µL PCR with primer sets that detect integrated DGR-ICE junction fragments or episomes
949 and chromosomal scars resulting from ICE excision (Figure 3A). To quantify ICE excision *in vivo*, germ-
950 free Swiss Webster mice (n=3) were monocolonized by *Bfr* carrying the plasmid pFD340, which confers
951 Erm resistance. Bacterial DNA was extracted from mouse feces and cecal content using the ZR fecal
952 DNA miniprep kit (Zymo Research). DNA from *Bfr* cells in scraped colon mucus was prepared using

953 DNeasy Blood and Tissue Kits (Qiagen) following pre-treatment with N-acetyl-L-cysteine. Specifically, a
954 freshly made NALC solution (50ml 2.94% sodium citrate, 50ml 4% sodium hydroxide, 500mg N-acetyl-
955 L-cysteine) was added to the mucus at 1:1 ratio (vol/vol), incubated at the room temperature for 1 hr with
956 agitation until the sample attained desired fluidity. Quantitative PCR (qPCR) was performed on an iCycler
957 iQ real-time PCR detection system (BioRad) with iQ SYBR Green Supermix (BioRad). Per each 30µl
958 qPCR, 30 ng DNA from various samples was used as template. Ct values for the episome and
959 chromosomal scar were normalized to Ct values of the housekeeping gene, rpoD (BF638R_RS13245).
960 Relative quantification of excised ICE in different samples was calculated by the $\Delta\Delta$ Ct approach, using
961 cultured *Bfr* with pFD340 as baseline. DNA extracted from *in vitro* grown cultures of *Bfr* carrying *paraC2*
962 was included as a positive control for high level excision.

963 **ICE Transfer assays.** To create ICEs and strains with compatible antibiotic markers, the Erm resistance
964 cassette in pKnock-erm was first swapped with a *tetQ* marker, resulting pKnock-tetQ. The DGR-ICE from
965 *Bfr* or *Bth* was tagged with the Tet resistance marker by inserting pKnock-tetQ downstream of *rt*. Mating
966 experiments were designed to measure inter- and intra- species ICE transfer. Briefly, three independent
967 cultures of donor or recipient cells were grown to mid-log phase (OD₆₀₀ 0.5), mixed at a 2:1 ratio and
968 spotted on sterile nitrocellulose membranes (0.45 µM, PALL Life Science). After incubation on non-
969 selective BHIS plates for 16-24 hrs, mating mixtures were washed off the filter into 2 mL of BHIS and
970 plated onto selective BHIS medium with 2 µg/mL Tet and 200 µg/mL FudR. Putative transconjugants
971 were purified and verified by PCR reactions with primer sets that are specific to the recipient or donor
972 ICE. The transfer efficiency was calculated by dividing the number of genuine transconjugants by the
973 number of recipients in each mating experiment. A similar transfer was also set up in gnotobiotic mice,
974 except that both the tetQ-tagged donor and Δ *tdk*ΔICE recipient carry pFD340 to confer erythromycin
975 resistance. Each mating pair was inoculated to 4 mice co-housed in one cage. Mouse feces were plated
976 at days 1, 3, 5, 7 and 10 post-inoculation.

977 **Assays for mutagenic retrohoming.** For *in vitro* assays, overnight cultures of *Bacteroides* strains were
978 diluted to OD₆₀₀ 0.01 in 3 mL of BHIS in triplicate. Every eight hours, the OD₆₀₀ was measured (OD₆₀₀
979 between 0.5-0.8), and cultures were rediluted to OD₆₀₀ 0.01 in 3 mL of fresh BHIS. Samples were
980 collected by pelleting 1 OD₆₀₀ of cells on Days 3, 7, 10, and 14 for Amplicon-Seq. *In vivo* assays were
981 conducted by collecting fecal pellets from monocolonized SW mice on Days 3, 7, 10, and 14 post-gavage.
982 Genomic DNA was extracted from pellets for Amplicon-Seq.

983 **Amplicon-Seq of VR regions.** Total DNA was extracted from bacterial cultures (PureLink Genomic DNA
984 Mini Kit, Thermo Fisher) or stool (QIAamp Fast DNA Stool Mini Kit, Qiagen). Primers amplifying the VR

985 region of each *Bacteroides* strain were designed with the forward primer containing a 20 bp random
986 Unique Molecular Index (UMI) and an adapter, as described previously¹¹¹. Briefly, each forward primer
987 was designed to anneal to the VR region at 62-64°C and was present in 1/10 the normal concentration.
988 A second forward primer would anneal to the adapter of the first forward primer at a temperature of 68-
989 70°C. The second forward primer and reverse primer contained partial Illumina adapter sequences.
990 Cycling parameters included 1 cycle of 98°C for 3 mins, 2 cycles of 98°C for 15 secs, 62°C for 45 secs,
991 and 72°C for 30 secs, followed by 38 cycles of 98°C for 15 secs, 70°C for 45 secs, followed by 72°C for
992 5 mins. PCR products were purified using SPRI Select beads (Beckman Coulter) and sent for EZ
993 Amplicon-Seq (Azenta Life Sciences).

994 **Data processing of Amplicon-Seq reads.** Raw reads were trimmed and contamination filtered with
995 BBDuk using default parameters except a kmer length of 23 and mink length of 11. Data was then merged
996 with BBMerge¹¹² using default parameters. Merged reads were then aligned to VR by creating a custom
997 blastn⁹⁹ database created from the parental VR sequences of *Bfr*, *Bth*, *Bun*, *Bov*, and *Bfi* to allow for a
998 large number of potential mismatches using the parameters for blastn word_size=8, reward=1, penalty=-
999 1, evalue=1e-5, gapopen=6, gapextend=6, and perc_identity=50. UMI and VR sequences were then
1000 extracted from the aligned read and compiled together. UMIs that contained mismatching sequences
1001 were discarded. Reads were then analyzed for the number of adenine-mutations and non-adenine
1002 mutations compared to parental strains. Reads that had >50% of non-adenine mutations were also
1003 discarded. The number of divergent reads was calculated by summing the number of reads with greater
1004 than 1 (*i.e.* 2 or more) adenine-mutations compared to parental strain (to account for potential sequencing
1005 errors). When plotted on a log chart, a pseudocount of 1 mutated read was added to all samples.

1006 **Total RNA extraction, RNA-Seq Library Preparation, and Illumina Sequencing.** *Bacteroides* cultures
1007 were grown to mid-log phase (OD₆₀₀ 0.3–0.5) and stationary phase (OD₆₀₀ 0.7–1.0) in BHIS. Cells were
1008 flash-frozen in a dry ice-ethanol slurry, then pelleted by centrifugation at 16,000xg for 1 minute at 4°C.
1009 The bacterial pellets were stored at –80°C until RNA extraction. RNA was extracted using TRIzol reagent
1010 and the PureLink™ RNA Mini Kit following the manufacturer's instructions. Briefly, bacterial pellets were
1011 incubated with 1 mL of TRIzol for 5 minutes at room temperature, followed by the addition of 200 µL of
1012 chloroform. After 5 minutes at room temperature, the samples were centrifuged at 12,000xg for 15
1013 minutes at 4°C. The aqueous phase was transferred to a new tube, and an equal volume of 100% ethanol
1014 was added. This mixture was transferred to a spin cartridge and centrifuged at 12,000xg for 15 seconds
1015 at room temperature.

1016 RNA samples were treated with PureLink™ DNase I, according to the manufacturer's protocol. Following
1017 DNase treatment and four wash steps, RNA was eluted in 50 µL of RNase-free water by centrifugation
1018 at 12,000xg for 1 minute. Eluted RNA was transferred to RNase-free tubes and analyzed for quality by
1019 agarose gel electrophoresis. RNA integrity was assessed using the Bio-Rad ChemiDoc system, and RNA
1020 concentration was quantified with a NanoDrop One spectrophotometer (Thermo Fisher Scientific).
1021 Further RNA quality and concentration was validated on an Agilent 4200 TapeStation RNA ScreenTape,
1022 with all RIN scores above 8.0 and samples normalized to 250 ng starting mass. Library preparation was
1023 performed using Illumina's Stranded Total RNA Prep with Ribo-Zero Plus Microbiome, which targets
1024 depletion of microbiome rRNA prior to cDNA synthesis and library formation. Final libraries were validated
1025 using an Agilent 4200 TapeStation D1000 ScreenTape, with an average fragment size of 396bp. Libraries
1026 were quantified using Invitrogen's Quant iT High Sensitivity dsDNA Kit and normalized during pooling.
1027 Samples were sequenced in 3 lanes of NovaSeq X Plus 10B PE 2x100.

1028 ***Bacteroides* dataset and Mother-Infant dataset acquisition.** The NCBI Assembly database was
1029 searched for the term "Bacteroides" and filtered to include those assemblies within the RefSeq database
1030 on 01/22/2021. Genome strain names were derived from the genome assembly reports. For the mother-
1031 infant dataset, raw reads were downloaded from BioProject PRJNA475246⁸⁵ or metagenomic
1032 assemblies^{84,96,97} were downloaded as previously described¹¹³. Raw reads had adapters trimmed and
1033 were quality filtered using BBDuk and merged using BBMerge¹¹². Merged reads were assembled into
1034 contigs using MegaHit²⁹ using parameters: --min-contig-len 500 -m 0.85 --presets meta-sensitive. Contigs
1035 with length <2000 bp were then discarded and the resulting contigs were used for DGR identification.
1036 Taxonomic classification of contigs was performed by mmseqs2 using default parameters¹¹⁴.

1037 **DGR identification from genomic and metagenomic datasets.** A profile of previously identified DGR
1038 RT proteins^{12,115} was built using HMMER⁵⁶. Predicted proteomes were derived from genomes and contigs
1039 using Prodigal⁹⁸ and the resulting predicted proteins were searched for hits using HMMScan. A 20 kb
1040 window on either side of RT hits was used as input to search for imperfect repeats using a custom BLAST
1041 script⁹⁹. Pairs were checked for one pair to be contained within a previously predicted ORF. Because
1042 there are many potential mismatches of reads to the VR sequence, traditional local aligners, such as
1043 Bowtie¹¹⁶ or BWA¹¹⁷ will trim or discard the read. To circumvent this challenge and determine DGR activity
1044 from metagenomic data, raw reads were mapped back to identified VRs by creating a custom BLAST
1045 database of the VR sequence using the same blastn parameters as above. Reads that fully or partially
1046 aligned were then checked for alignment to TR and the rest of the genome/contig by identifying if those
1047 reads had fewer mismatches than to VR. Any read which best aligned to the VR region was kept. Reads
1048 aligning to VR regions were then analyzed to determine mutations existed that corresponded to TR

1049 adenine positions. If adenine-specific mutations were found, the DGR was categorized as active. VR
1050 haplotypes were generated for each timepoint by generating the consensus VR sequence. The
1051 consensus sequence was compared between timepoints to determine if it differed.

1052 **Identification of non-DGR RTs.** A previously built profile of all RTs in bacterial genomes⁵¹ was modified
1053 to exclude the DGR RTs and rebuilt using HMMER. Proteomes from genomes or metagenomic contigs
1054 were searched for these RTs. The number of non-DGR RTs was reduced by random sampling and this
1055 subset was added to the identified DGR RTs for phylogenetic analysis.

1056 **Clustering of variable proteins and phylogenetic tree building.** In order to cluster variable proteins,
1057 each primary amino acid sequence was used for an all-vs-all blastp⁹⁹ search. The result of the blastp was
1058 used as an edge weight for input into MCL clustering with inflation value of 2.0¹⁰¹. Within each cluster, a
1059 multiple sequence alignment was performed using MUSCLE¹¹⁸. This alignment was input into HHpred⁵⁴
1060 to search for functional domains within the proteins using the databases PDB_mmCIF70_17_Apr, Pfam-
1061 A_v35, and COG-KOG_v1.0. The RT phylogenetic tree was created by aligning each RT protein
1062 sequence, which was input into Fasttree2 using default parameters¹⁰⁰. Variable proteins were classified
1063 as pilus proteins if they shared a significant domain with the following PDBs: 4EPS⁴⁸, 4QB7⁴⁸, 6JZJ⁶³,
1064 5NF4¹¹⁹. Variable proteins were classified as TaqVP-like if they shared a significant domain with PDB:
1065 5VF4⁵⁷. Variable proteins were classified as phage receptor tail binding proteins if they shared a
1066 significant domain with PDB: 1YU0¹⁵.

1067 **DGR genomic assignments.** Contigs containing DGRs were used as input for geNomad⁵², which
1068 predicts regions of contigs to be viral or plasmid. Coordinates from those predictions were aligned to
1069 DGR coordinates. DGRs located on a contig predicted to be entirely viral were classified as virus. DGRs
1070 located within a viral region surrounded by bacterial genes were classified as prophage. A 100 kb window
1071 surrounding DGRs located within predicted plasmid regions was then used as input for ICEfinder⁵³ for
1072 ICE prediction. ICE synteny figures were generated using pyGenomeViz¹²⁰.

1073 **DGR gene transcription analysis.** Raw sequencing reads from RNA-seq experiments were quality
1074 filtered and trimmed with fastp v0.23.4¹²¹ with default parameters, then mapped to the reference genome
1075 of each respective species with bwa-mem2¹²² with default parameters. Mapping statistics were generated
1076 with bamtools¹²³ v2.5.1 and visualized with multiqc¹²⁴ v1.21. Reads in gene features were counted with
1077 featureCounts¹²⁵ function of the subread package at the fragment level (--countReadPairs) with
1078 fragments overlapping multiple features counted in both (-O). Total transcript counts were normalized to
1079 trimmed mean of M (TMM) values using NOISeq¹²⁶ with default parameters. Within each *Bacteroides*
1080 strain, DGR gene TMM values were normalized to *gyrA* gene TMM values. The ratio of DGR genes to

1081 *gyrA* were then compared across strains, setting *Bfr* ratios to 1 to generate relative ratios for each of the
1082 other strains.

1083 **VR mutagenesis and entropy analysis.** VR mutagenesis frequency was calculated by taking the
1084 number of reads with >1 substitution from the canonical sequence divided by the total number of
1085 sequences. VR entropy was calculated using the formula:

1086
$$Entropy = \sum_1^n p * \log(p)$$

1087 according to Shannon⁸¹, where p is the frequency of a unique VR sequence and n is the total number of
1088 unique mutagenized VR sequences.

1089 **Statistical tests and metrics.** Statistical comparisons were performed using the tests indicated in the
1090 figure legends. P-values were generated using ANOVA or student's t-test. Multiple comparison
1091 corrections were performed using Holm-Sidak method or Tukey's method, where appropriate. Statistical
1092 significance was defined at $\alpha=0.05$.

1093 **Data and code availability.** Python code to identify DGRs from genomes will be made available at
1094 <https://github.com/macadangdanglab/dgrdiscovery> upon publication. Amplicon-Seq reads of VRs will be
1095 accessible on the short read archive (SRA) upon publication.

1096 **References**

1097

1098 1. Heilbron, K., Toll-Riera, M., Kojadinovic, M., and MacLean, R.C. (2014). Fitness Is Strongly
1099 Influenced by Rare Mutations of Large Effect in a Microbial Mutation Accumulation Experiment.
1100 *Genetics* 197, 981–990. <https://doi.org/10.1534/genetics.114.163147>.

1101 2. Mazel, D. (2006). Integrons: agents of bacterial evolution. *Nat Rev Microbiol* 4, 608–620.
1102 <https://doi.org/10.1038/nrmicro1462>.

1103 3. Müller, F., and Tobler, H. (2000). Chromatin diminution in the parasitic nematodes *Ascaris suum* and
1104 *Parascaris univalens*. *Int J Parasitol* 30, 391–399. [https://doi.org/10.1016/s0020-7519\(99\)00199-x](https://doi.org/10.1016/s0020-7519(99)00199-x).

1105 4. Schatz, D.G., and Swanson, P.C. (2011). V(D)J Recombination: Mechanisms of Initiation. *Annu Rev*
1106 *Genet* 45, 167–202. <https://doi.org/10.1146/annurev-genet-110410-132552>.

1107 5. Fitzgerald, D.M., and Rosenberg, S.M. (2019). What is mutation? A chapter in the series: How
1108 microbes “jeopardize” the modern synthesis. *Plos Genet* 15, e1007995.
1109 <https://doi.org/10.1371/journal.pgen.1007995>.

1110 6. Shee, C., Gibson, J.L., Darrow, M.C., Gonzalez, C., and Rosenberg, S.M. (2011). Impact of a stress-
1111 inducible switch to mutagenic repair of DNA breaks on mutation in *Escherichia coli*. *Proc National Acad*
1112 *Sci* 108, 13659–13664. <https://doi.org/10.1073/pnas.1104681108>.

1113 7. Jiang, F., and Doudna, J.A. (2015). CRISPR–Cas9 Structures and Mechanisms. *Annu. Rev.*
1114 *Biophys.* 46, 1–25. <https://doi.org/10.1146/annurev-biophys-062215-010822>.

1115 8. Marraffini, L.A., and Sontheimer, E.J. (2008). CRISPR Interference Limits Horizontal Gene Transfer
1116 in *Staphylococci* by Targeting DNA. *Science* 322, 1843–1845. <https://doi.org/10.1126/science.1165771>.

1117 9. Wheatley, R.M., and MacLean, R.C. (2021). CRISPR-Cas systems restrict horizontal gene transfer in
1118 *Pseudomonas aeruginosa*. *ISME J.* 15, 1420–1433. <https://doi.org/10.1038/s41396-020-00860-3>.

1119 10. Macadangang, B.R., Makanani, S.K., and Miller, J.F. (2022). Accelerated Evolution by Diversity-
1120 Generating Retroelements. *Annu Rev Microbiol* 76, 389–411. <https://doi.org/10.1146/annurev-micro-030322-040423>.

1121

- 1122 11. Doré, H., Eisenberg, A.R., Junkins, E.N., Leventhal, G.E., Ganesh, A., Cordero, O.X., Paul, B.G.,
1123 Valentine, D.L., O'Malley, M.A., and Wilbanks, E.G. (2024). Targeted hypermutation of putative antigen
1124 sensors in multicellular bacteria. *Proc. Natl. Acad. Sci.* *121*, e2316469121.
1125 <https://doi.org/10.1073/pnas.2316469121>.
- 1126 12. Wu, L., Gingery, M., Abebe, M., Arambula, D., Czornyj, E., Handa, S., Khan, H., Liu, M.,
1127 Pohlschroder, M., Shaw, K.L., et al. (2017). Diversity-generating retroelements: natural variation,
1128 classification and evolution inferred from a large-scale genomic survey. *Nucleic Acids Res* *46*, gkx1150-
1129 . <https://doi.org/10.1093/nar/gkx1150>.
- 1130 13. Miller, J.L., Coq, J.L., Hodes, A., Barbalat, R., Miller, J.F., and Ghosh, P. (2008). Selective Ligand
1131 Recognition by a Diversity-Generating Retroelement Variable Protein. *Plos Biol* *6*, e131.
1132 <https://doi.org/10.1371/journal.pbio.0060131>.
- 1133 14. Coq, J.L., and Ghosh, P. (2011). Conservation of the C-type lectin fold for massive sequence
1134 variation in a *Treponema* diversity-generating retroelement. *Proc National Acad Sci* *108*, 14649–14653.
1135 <https://doi.org/10.1073/pnas.1105613108>.
- 1136 15. McMahon, S.A., Miller, J.L., Lawton, J.A., Kerkow, D.E., Hodes, A., Marti-Renom, M.A., Doulatov,
1137 S., Narayanan, E., Sali, A., Miller, J.F., et al. (2005). The C-type lectin fold as an evolutionary solution
1138 for massive sequence variation. *Nat Struct Mol Biol* *12*, 886–892. <https://doi.org/10.1038/nsmb992>.
- 1139 16. Guo, H., Tse, L.V., Barbalat, R., Sivaamnuaiphorn, S., Xu, M., Doulatov, S., and Miller, J.F. (2008).
1140 Diversity-Generating Retroelement Homing Regenerates Target Sequences for Repeated Rounds of
1141 Codon Rewriting and Protein Diversification. *Mol Cell* *31*, 813–823.
1142 <https://doi.org/10.1016/j.molcel.2008.07.022>.
- 1143 17. Doulatov, S., Hodes, A., Dai, L., Mandhana, N., Liu, M., Deora, R., Simons, R.W., Zimmerly, S., and
1144 Miller, J.F. (2004). Tropism switching in *Bordetella* bacteriophage defines a family of diversity-
1145 generating retroelements. *Nature* *431*, 476–481. <https://doi.org/10.1038/nature02833>.
- 1146 18. Liu, M., Deora, R., Doulatov, S.R., Gingery, M., Eiserling, F.A., Preston, A., Maskell, D.J., Simons,
1147 R.W., Cotter, P.A., Parkhill, J., et al. (2002). Reverse Transcriptase-Mediated Tropism Switching in
1148 *Bordetella* Bacteriophage. *Science* *295*, 2091–2094. <https://doi.org/10.1126/science.1067467>.

- 1149 19. Roux, S., Paul, B.G., Bagby, S.C., Nayfach, S., Allen, M.A., Attwood, G., Cavicchioli, R.,
1150 Chistoserdova, L., Gruninger, R.J., Hallam, S.J., et al. (2021). Ecology and molecular targets of
1151 hypermutation in the global microbiome. *Nat Commun* 12, 3076. [https://doi.org/10.1038/s41467-021-](https://doi.org/10.1038/s41467-021-23402-7)
1152 [23402-7](https://doi.org/10.1038/s41467-021-23402-7).
- 1153 20. Macadangdang, B.R., Makanani, S.K., and Miller, J.F. (2022). Accelerated Evolution by Diversity-
1154 Generating Retroelements. *Annu Rev Microbiol* 76. [https://doi.org/10.1146/annurev-micro-030322-](https://doi.org/10.1146/annurev-micro-030322-040423)
1155 [040423](https://doi.org/10.1146/annurev-micro-030322-040423).
- 1156 21. Faith, J.J., Guruge, J.L., Charbonneau, M., Subramanian, S., Seedorf, H., Goodman, A.L.,
1157 Clemente, J.C., Knight, R., Heath, A.C., Leibel, R.L., et al. (2013). The Long-Term Stability of the
1158 Human Gut Microbiota. *Science* 341, 1237439. <https://doi.org/10.1126/science.1237439>.
- 1159 22. Charbonneau, M.R., O'Donnell, D., Blanton, L.V., Totten, S.M., Davis, J.C.C., Barratt, M.J., Cheng,
1160 J., Guruge, J., Talcott, M., Bain, J.R., et al. (2016). Sialylated Milk Oligosaccharides Promote
1161 Microbiota-Dependent Growth in Models of Infant Undernutrition. *Cell* 164, 859–871.
1162 <https://doi.org/10.1016/j.cell.2016.01.024>.
- 1163 23. Mazmanian, S.K., Round, J.L., and Kasper, D.L. (2008). A microbial symbiosis factor prevents
1164 intestinal inflammatory disease. *Nature* 453, 620–625. <https://doi.org/10.1038/nature07008>.
- 1165 24. Portincasa, P., Bonfrate, L., Vacca, M., Angelis, M.D., Farella, I., Lanza, E., Khalil, M., Wang, D.Q.-
1166 H., Sperandio, M., and Ciaula, A.D. (2022). Gut Microbiota and Short Chain Fatty Acids: Implications in
1167 Glucose Homeostasis. *Int. J. Mol. Sci.* 23, 1105. <https://doi.org/10.3390/ijms23031105>.
- 1168 25. Jean, S., Wallace, M.J., Dantas, G., and Burnham, C.-A.D. (2022). Time for Some Group Therapy:
1169 Update on Identification, Antimicrobial Resistance, Taxonomy, and Clinical Significance of the
1170 *Bacteroides fragilis* Group. *J. Clin. Microbiol.* 60, e02361-20. <https://doi.org/10.1128/jcm.02361-20>.
- 1171 26. Schloissnig, S., Arumugam, M., Sunagawa, S., Mitreva, M., Tap, J., Zhu, A., Waller, A., Mende,
1172 D.R., Kultima, J.R., Martin, J., et al. (2013). Genomic variation landscape of the human gut microbiome.
1173 *Nature* 493, 45–50. <https://doi.org/10.1038/nature11711>.
- 1174 27. Liebert, C.A., Hall, R.M., and Summers, A.O. (1999). Transposon Tn 21 , Flagship of the Floating
1175 Genome. *Microbiol Mol Biol R* 63, 507–522. <https://doi.org/10.1128/mubr.63.3.507-522.1999>.

- 1176 28. Frost, L.S., Leplae, R., Summers, A.O., and Toussaint, A. (2005). Mobile genetic elements: the
1177 agents of open source evolution. *Nat Rev Microbiol* 3, 722–732. <https://doi.org/10.1038/nrmicro1235>.
- 1178 29. Li, D., Liu, C.-M., Luo, R., Sadakane, K., and Lam, T.-W. (2015). MEGAHIT: an ultra-fast single-
1179 node solution for large and complex metagenomics assembly via succinct de Bruijn graph.
1180 *Bioinformatics* 31, 1674–1676. <https://doi.org/10.1093/bioinformatics/btv033>.
- 1181 30. Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S., Lesin, V.M.,
1182 Nikolenko, S.I., Pham, S., Prjibelski, A.D., et al. (2012). SPAdes: A New Genome Assembly Algorithm
1183 and Its Applications to Single-Cell Sequencing. *J Comput Biol* 19, 455–477.
1184 <https://doi.org/10.1089/cmb.2012.0021>.
- 1185 31. Nurk, S., Meleshko, D., Korobeynikov, A., and Pevzner, P.A. (2017). metaSPAdes: a new versatile
1186 metagenomic assembler. *Genome Res* 27, 824–834. <https://doi.org/10.1101/gr.213959.116>.
- 1187 32. Boisvert, S., Raymond, F., Godzaridis, É., Laviolette, F., and Corbeil, J. (2012). Ray Meta: scalable
1188 de novo metagenome assembly and profiling. *Genome Biol* 13, R122. <https://doi.org/10.1186/gb-2012-13-12-r122>.
- 1190 33. Afiahayati, Sato, K., and Sakakibara, Y. (2015). MetaVelvet-SL: an extension of the Velvet
1191 assembler to a de novo metagenomic assembler utilizing supervised learning. *Dna Res* 22, 69–77.
1192 <https://doi.org/10.1093/dnares/dsu041>.
- 1193 34. Namiki, T., Hachiya, T., Tanaka, H., and Sakakibara, Y. (2012). MetaVelvet: an extension of Velvet
1194 assembler to de novo metagenome assembly from short sequence reads. *Nucleic Acids Res* 40, e155–
1195 e155. <https://doi.org/10.1093/nar/gks678>.
- 1196 35. Peng, Y., Leung, H.C.M., Yiu, S.M., and Chin, F.Y.L. (2012). IDBA-UD: a de novo assembler for
1197 single-cell and metagenomic sequencing data with highly uneven depth. *Bioinformatics* 28, 1420–1428.
1198 <https://doi.org/10.1093/bioinformatics/bts174>.
- 1199 36. Treangen, T.J., Koren, S., Sommer, D.D., Liu, B., Astrovskaya, I., Ondov, B., Darling, A.E.,
1200 Phillippy, A.M., and Pop, M. (2013). MetAMOS: a modular and open source metagenomic assembly
1201 and analysis pipeline. *Genome Biol* 14, R2. <https://doi.org/10.1186/gb-2013-14-1-r2>.

- 1202 37. Kultima, J.R., Coelho, L.P., Forslund, K., Huerta-Cepas, J., Li, S.S., Driessen, M., Voigt, A.Y.,
1203 Zeller, G., Sunagawa, S., and Bork, P. (2016). MOCAT2: a metagenomic assembly, annotation and
1204 profiling framework. *Bioinformatics* 32, 2520–2523. <https://doi.org/10.1093/bioinformatics/btw183>.
- 1205 38. Eren, A.M., Esen, Ö.C., Quince, C., Vineis, J.H., Morrison, H.G., Sogin, M.L., and Delmont, T.O.
1206 (2015). Anvi'o: an advanced analysis and visualization platform for 'omics data. *Peerj* 3, e1319.
1207 <https://doi.org/10.7717/peerj.1319>.
- 1208 39. Wu, Y.-W., Simmons, B.A., and Singer, S.W. (2016). MaxBin 2.0: an automated binning algorithm
1209 to recover genomes from multiple metagenomic datasets. *Bioinformatics* 32, 605–607.
1210 <https://doi.org/10.1093/bioinformatics/btv638>.
- 1211 40. Alneberg, J., Bjarnason, B.S., Bruijn, I. de, Schirmer, M., Quick, J., Ijaz, U.Z., Lahti, L., Loman, N.J.,
1212 Andersson, A.F., and Quince, C. (2014). Binning metagenomic contigs by coverage and composition.
1213 *Nat Methods* 11, 1144–1146. <https://doi.org/10.1038/nmeth.3103>.
- 1214 41. Lu, Y.Y., Chen, T., Fuhrman, J.A., and Sun, F. (2016). COCACOLA: binning metagenomic contigs
1215 using sequence COmposition, read CoverAge, CO-alignment and paired-end read LinkAge.
1216 *Bioinformatics* 33, btw290. <https://doi.org/10.1093/bioinformatics/btw290>.
- 1217 42. Kang, D.D., Froula, J., Egan, R., and Wang, Z. (2015). MetaBAT, an efficient tool for accurately
1218 reconstructing single genomes from complex microbial communities. *Peerj* 3, e1165.
1219 <https://doi.org/10.7717/peerj.1165>.
- 1220 43. Laczny, C.C., Sternal, T., Plugaru, V., Gawron, P., Atashpendar, A., Margossian, H.H., Coronado,
1221 S., Maaten, L. van der, Vlassis, N., and Wilmes, P. (2015). VizBin - an application for reference-
1222 independent visualization and human-augmented binning of metagenomic data. *Microbiome* 3, 1.
1223 <https://doi.org/10.1186/s40168-014-0066-1>.
- 1224 44. Wu, Y.-W., and Ye, Y. (2011). A Novel Abundance-Based Algorithm for Binning Metagenomic
1225 Sequences Using I-tuples. *J Comput Biol* 18, 523–534. <https://doi.org/10.1089/cmb.2010.0245>.
- 1226 45. Imelfort, M., Parks, D., Woodcroft, B.J., Dennis, P., Hugenholtz, P., and Tyson, G.W. (2014).
1227 GroopM: an automated tool for the recovery of population genomes from related metagenomes. *Peerj*
1228 2, e603. <https://doi.org/10.7717/peerj.603>.

- 1229 46. Wang, Y., Leung, H.C.M., Yiu, S.M., and Chin, F.Y.L. (2012). MetaCluster 5.0: a two-round binning
1230 approach for metagenomic data for low-abundance species in a noisy sample. *Bioinformatics* 28, i356–
1231 i362. <https://doi.org/10.1093/bioinformatics/bts397>.
- 1232 47. Patil, K.R., Roune, L., and McHardy, A.C. (2012). The PhyloPythiaS Web Server for Taxonomic
1233 Assignment of Metagenome Sequences. *Plos One* 7, e38581.
1234 <https://doi.org/10.1371/journal.pone.0038581>.
- 1235 48. Xu, Q., Shoji, M., Shibata, S., Naito, M., Sato, K., Elsliger, M.-A., Grant, J.C., Axelrod, H.L., Chiu,
1236 H.-J., Farr, C.L., et al. (2016). A Distinct Type of Pilus from the Human Microbiome. *Cell* 165, 690–703.
1237 <https://doi.org/10.1016/j.cell.2016.03.016>.
- 1238 49. O’Leary, N.A., Wright, M.W., Brister, J.R., Ciuffo, S., Haddad, D., McVeigh, R., Rajput, B.,
1239 Robbertse, B., Smith-White, B., Ako-Adjei, D., et al. (2016). Reference sequence (RefSeq) database at
1240 NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* 44, D733–
1241 D745. <https://doi.org/10.1093/nar/gkv1189>.
- 1242 50. Vallota-Eastman, A., Arrington, E.C., Meeken, S., Roux, S., Dasari, K., Rosen, S., Miller, J.F.,
1243 Valentine, D.L., and Paul, B.G. (2020). Role of diversity-generating retroelements for regulatory
1244 pathway tuning in cyanobacteria. *Bmc Genomics* 21, 664. <https://doi.org/10.1186/s12864-020-07052-5>.
- 1245 51. Sharifi, F., and Ye, Y. (2021). Identification and classification of reverse transcriptases in bacterial
1246 genomes and metagenomes. *Nucleic Acids Res* 50, e29–e29. <https://doi.org/10.1093/nar/gkab1207>.
- 1247 52. Camargo, A.P., Roux, S., Schulz, F., Babinski, M., Xu, Y., Hu, B., Chain, P.S.G., Nayfach, S., and
1248 Kyrpides, N.C. (2023). Identification of mobile genetic elements with geNomad. *Nat. Biotechnol.*, 1–10.
1249 <https://doi.org/10.1038/s41587-023-01953-y>.
- 1250 53. Liu, M., Li, X., Xie, Y., Bi, D., Sun, J., Li, J., Tai, C., Deng, Z., and Ou, H.-Y. (2019). ICEberg 2.0: an
1251 updated database of bacterial integrative and conjugative elements. *Nucleic Acids Res* 47, D660–
1252 D665. <https://doi.org/10.1093/nar/gky1123>.
- 1253 54. Söding, J., Biegert, A., and Lupas, A.N. (2005). The HHpred interactive server for protein homology
1254 detection and structure prediction. *Nucleic Acids Res* 33, W244–W248.
1255 <https://doi.org/10.1093/nar/gki408>.

- 1256 55. Paul, B.G., Bagby, S.C., Czornyj, E., Arambula, D., Handa, S., Sczyrba, A., Ghosh, P., Miller, J.F.,
1257 and Valentine, D.L. (2015). Targeted diversity generation by intraterrestrial archaea and archaeal
1258 viruses. *Nat Commun* 6, 6585. <https://doi.org/10.1038/ncomms7585>.
- 1259 56. Eddy, S.R. (2011). Accelerated Profile HMM Searches. *Plos Comput Biol* 7, e1002195.
1260 <https://doi.org/10.1371/journal.pcbi.1002195>.
- 1261 57. Handa, S., Shaw, K.L., and Ghosh, P. (2019). Crystal structure of a *Thermus aquaticus* diversity-
1262 generating retroelement variable protein. *Plos One* 14, e0205618.
1263 <https://doi.org/10.1371/journal.pone.0205618>.
- 1264 58. Nguyen, K.B., Sreelatha, A., Durrant, E.S., Lopez-Garrido, J., Muszewska, A., Dudkiewicz, M.,
1265 Grynberg, M., Yee, S., Pogliano, K., Tomchick, D.R., et al. (2016). Phosphorylation of spore coat
1266 proteins by a family of atypical protein kinases. *Proc National Acad Sci* 113, E3482–E3491.
1267 <https://doi.org/10.1073/pnas.1605917113>.
- 1268 59. Welch, J.L.M., Rossetti, B.J., Rieken, C.W., Dewhirst, F.E., and Borisy, G.G. (2016). Biogeography
1269 of a human oral microbiome at the micron scale. *Proc National Acad Sci* 113, E791–E800.
1270 <https://doi.org/10.1073/pnas.1522149113>.
- 1271 60. Welch, J.L.M., Ramírez-Puebla, S.T., and Borisy, G.G. (2020). Oral Microbiome Geography:
1272 Micron-Scale Habitat and Niche. *Cell Host Microbe* 28, 160–168.
1273 <https://doi.org/10.1016/j.chom.2020.07.009>.
- 1274 61. Alayyoubi, M., Guo, H., Dey, S., Golnazarian, T., Brooks, G.A., Rong, A., Miller, J.F., and Ghosh, P.
1275 (2013). Structure of the Essential Diversity-Generating Retroelement Protein bAvd and Its Functionally
1276 Important Interaction with Reverse Transcriptase. *Structure* 21, 266–276.
1277 <https://doi.org/10.1016/j.str.2012.11.016>.
- 1278 62. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K.,
1279 Bates, R., Žídek, A., Potapenko, A., et al. (2021). Highly accurate protein structure prediction with
1280 AlphaFold. *Nature* 596, 583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- 1281 63. Shibata, S., Shoji, M., Okada, K., Matsunami, H., Matthews, M.M., Imada, K., Nakayama, K., and
1282 Wolf, M. (2020). Structure of polymerized type V pilin reveals assembly mechanism involving protease-
1283 mediated strand exchange. *Nat Microbiol* 5, 830–837. <https://doi.org/10.1038/s41564-020-0705-1>.

- 1284 64. Hospenthal, M.K., Costa, T.R.D., and Waksman, G. (2017). A comprehensive guide to pilus
1285 biogenesis in Gram-negative bacteria. *Nat Rev Microbiol* 15, 365–379.
1286 <https://doi.org/10.1038/nrmicro.2017.40>.
- 1287 65. Götzke, H., Kilisch, M., Martínez-Carranza, M., Sograte-Idrissi, S., Rajavel, A., Schlichthaerle, T.,
1288 Engels, N., Jungmann, R., Stenmark, P., Opazo, F., et al. (2019). The ALFA-tag is a highly versatile
1289 tool for nanobody-based bioscience applications. *Nat Commun* 10, 4403.
1290 <https://doi.org/10.1038/s41467-019-12301-7>.
- 1291 66. Johnson, C.M., and Grossman, A.D. (2015). Integrative and Conjugative Elements (ICEs): What
1292 They Do and How They Work. *Annu Rev Genet* 49, 1–25. [https://doi.org/10.1146/annurev-genet-](https://doi.org/10.1146/annurev-genet-112414-055018)
1293 [112414-055018](https://doi.org/10.1146/annurev-genet-112414-055018).
- 1294 67. Durrant, M.G., Li, M.M., Siranosian, B.A., Montgomery, S.B., and Bhatt, A.S. (2020). A
1295 Bioinformatic Analysis of Integrative Mobile Genetic Elements Highlights Their Role in Bacterial
1296 Adaptation. *Cell Host Microbe* 27, 140-153.e9. <https://doi.org/10.1016/j.chom.2019.10.022>.
- 1297 68. Franke, A.E., and Clewell, D.B. (1981). Evidence for a chromosome-borne resistance transposon
1298 (Tn916) in *Streptococcus faecalis* that is capable of “conjugal” transfer in the absence of a conjugative
1299 plasmid. *J. Bacteriol.* 145, 494–502. <https://doi.org/10.1128/jb.145.1.494-502.1981>.
- 1300 69. Mays, T.D., Smith, C.J., Welch, R.A., Delfini, C., and Macrina, F.L. (1982). Novel antibiotic
1301 resistance transfer in *Bacteroides*. *Antimicrob. Agents Chemother.* 21, 110–118.
1302 <https://doi.org/10.1128/aac.21.1.110>.
- 1303 70. Roberts, M.C., and Smith, A.L. (1980). Molecular characterization of “plasmid-free” antibiotic-
1304 resistant *Haemophilus influenzae*. *J. Bacteriol.* 144, 476–479. [https://doi.org/10.1128/jb.144.1.476-](https://doi.org/10.1128/jb.144.1.476-479.1980)
1305 [479.1980](https://doi.org/10.1128/jb.144.1.476-479.1980).
- 1306 71. Shoemaker, N.B., Smith, M.D., and Guild, W.R. (1980). DNase-resistant transfer of chromosomal
1307 cat and tet insertions by filter mating in pneumococcus. *Plasmid* 3, 80–87.
1308 [https://doi.org/10.1016/s0147-619x\(80\)90036-0](https://doi.org/10.1016/s0147-619x(80)90036-0).
- 1309 72. Botelho, J., and Schulenburg, H. (2021). The Role of Integrative and Conjugative Elements in
1310 Antibiotic Resistance Evolution. *Trends Microbiol.* 29, 8–18. <https://doi.org/10.1016/j.tim.2020.05.011>.

1311 73. Jaworski, D.D., and Clewell, D.B. (1994). Evidence that coupling sequences play a frequency-
1312 determining role in conjugative transposition of Tn916 in *Enterococcus faecalis*. *J. Bacteriol.* 176,
1313 3328–3335. <https://doi.org/10.1128/jb.176.11.3328-3335.1994>.

1314 74. Park, J., and Salyers, A.A. (2011). Characterization of the *Bacteroides* CTnDOT Regulatory Protein
1315 RteC. *J Bacteriol* 193, 91–97. <https://doi.org/10.1128/jb.01015-10>.

1316 75. Wexler, A.G., and Goodman, A.L. (2017). An insider's perspective: *Bacteroides* as a window into
1317 the microbiome. *Nat Microbiol* 2, 17026. <https://doi.org/10.1038/nmicrobiol.2017.26>.

1318 76. Russell, A.B., Wexler, A.G., Harding, B.N., Whitney, J.C., Bohn, A.J., Goo, Y.A., Tran, B.Q., Barry,
1319 N.A., Zheng, H., Peterson, S.B., et al. (2014). A Type VI Secretion-Related Pathway in *Bacteroidetes*
1320 Mediates Interbacterial Antagonism. *Cell Host Microbe* 16, 227–236.
1321 <https://doi.org/10.1016/j.chom.2014.07.007>.

1322 77. Arambula, D., Wong, W., Medhekar, B.A., Guo, H., Gingery, M., Czornyj, E., Liu, M., Dey, S.,
1323 Ghosh, P., and Miller, J.F. (2013). Surface display of a massively variable lipoprotein by a *Legionella*
1324 diversity-generating retroelement. *Proc National Acad Sci* 110, 8212–8217.
1325 <https://doi.org/10.1073/pnas.1301366110>.

1326 78. Naorem, S.S., Han, J., Wang, S., Lee, W.R., Heng, X., Miller, J.F., and Guo, H. (2017). DGR
1327 mutagenic transposition occurs via hypermutagenic reverse transcription primed by nicked template
1328 RNA. *Proc National Acad Sci* 114, E10187–E10195. <https://doi.org/10.1073/pnas.1715952114>.

1329 79. Handa, S., Reyna, A., Wiryaman, T., and Ghosh, P. (2020). Determinants of adenine-mutagenesis
1330 in diversity-generating retroelements. *Nucleic Acids Res.* 49, 1033–1045.
1331 <https://doi.org/10.1093/nar/gkaa1240>.

1332 80. Handa, S., Jiang, Y., Tao, S., Foreman, R., Schinazi, R.F., Miller, J.F., and Ghosh, P. (2018).
1333 Template-assisted synthesis of adenine-mutagenized cDNA by a retroelement protein complex. *Nucleic*
1334 *Acids Res* 46, gky620-. <https://doi.org/10.1093/nar/gky620>.

1335 81. Shannon, C.E. (1948). A mathematical theory of communication. *Bell Syst Technical J* 27, 379–
1336 423. <https://doi.org/10.1002/j.1538-7305.1948.tb01338.x>.

- 1337 82. Konopiński, M.K. (2020). Shannon diversity index: a call to replace the original Shannon's formula
1338 with unbiased estimator in the population genetics studies. PeerJ 8, e9391.
1339 <https://doi.org/10.7717/peerj.9391>.
- 1340 83. Brand, M.W., Wannemuehler, M.J., Phillips, G.J., Proctor, A., Overstreet, A.-M., Jergens, A.E.,
1341 Orcutt, R.P., and Fox, J.G. (2015). The Altered Schaedler Flora: Continued Applications of a Defined
1342 Murine Microbial Community. ILAR J. 56, 169–178. <https://doi.org/10.1093/ilar/ilv012>.
- 1343 84. Ferretti, P., Pasolli, E., Tett, A., Asnicar, F., Gorfer, V., Fedi, S., Armanini, F., Truong, D.T., Manara,
1344 S., Zolfo, M., et al. (2018). Mother-to-Infant Microbial Transmission from Different Body Sites Shapes
1345 the Developing Infant Gut Microbiome. Cell Host Microbe 24, 133-145.e5.
1346 <https://doi.org/10.1016/j.chom.2018.06.005>.
- 1347 85. Yassour, M., Jason, E., Hogstrom, L.J., Arthur, T.D., Tripathi, S., Siljander, H., Selvenius, J.,
1348 Oikarinen, S., Hyöty, H., Virtanen, S.M., et al. (2018). Strain-Level Analysis of Mother-to-Child Bacterial
1349 Transmission during the First Few Months of Life. Cell Host Microbe 24, 146-154.e4.
1350 <https://doi.org/10.1016/j.chom.2018.06.007>.
- 1351 86. Bäckhed, F., Ding, H., Wang, T., Hooper, L.V., Koh, G.Y., Nagy, A., Semenkovich, C.F., and
1352 Gordon, J.I. (2004). The gut microbiota as an environmental factor that regulates fat storage. Proc
1353 National Acad Sci 101, 15718–15723. <https://doi.org/10.1073/pnas.0407076101>.
- 1354 87. Lloyd-Price, J., Mahurkar, A., Rahnavard, G., Crabtree, J., Orvis, J., Hall, A.B., Brady, A., Creasy,
1355 H.H., McCracken, C., Giglio, M.G., et al. (2017). Strains, functions and dynamics in the expanded
1356 Human Microbiome Project. Nature 550, 61–66. <https://doi.org/10.1038/nature23889>.
- 1357 88. Zepeda-Rivera, M., Minot, S.S., Bouzek, H., Wu, H., Blanco-Míguez, A., Manghi, P., Jones, D.S.,
1358 LaCourse, K.D., Wu, Y., McMahon, E.F., et al. (2024). A distinct *Fusobacterium nucleatum* clade
1359 dominates the colorectal cancer niche. Nature, 1–9. <https://doi.org/10.1038/s41586-024-07182-w>.
- 1360 89. Yaffe, E., and Relman, D.A. (2020). Tracking microbial evolution in the human gut using Hi-C
1361 reveals extensive horizontal gene transfer, persistence and adaptation. Nat. Microbiol. 5, 343–353.
1362 <https://doi.org/10.1038/s41564-019-0625-0>.

- 1363 90. Zahavi, L., Lavon, A., Reicher, L., Shoer, S., Godneva, A., Leviatan, S., Rein, M., Weissbrod, O.,
1364 Weinberger, A., and Segal, E. (2023). Bacterial SNPs in the human gut microbiome associate with host
1365 BMI. *Nat. Med.* 29, 2785–2792. <https://doi.org/10.1038/s41591-023-02599-8>.
- 1366 91. Zhao, S., Lieberman, T.D., Poyet, M., Kauffman, K.M., Gibbons, S.M., Groussin, M., Xavier, R.J.,
1367 and Alm, E.J. (2019). Adaptive Evolution within Gut Microbiomes of Healthy People. *Cell Host Microbe*
1368 25, 656-667.e8. <https://doi.org/10.1016/j.chom.2019.03.007>.
- 1369 92. Garud, N.R., Good, B.H., Hallatschek, O., and Pollard, K.S. (2019). Evolutionary dynamics of
1370 bacteria in the gut microbiome within and across hosts. *PLoS Biol.* 17, e3000102.
1371 <https://doi.org/10.1371/journal.pbio.3000102>.
- 1372 93. Wolff, R., and Garud, N.R. (2023). Pervasive selective sweeps across human gut microbiomes.
1373 bioRxiv, 2023.12.22.573162. <https://doi.org/10.1101/2023.12.22.573162>.
- 1374 94. Valles-Colomer, M., Blanco-Míguez, A., Manghi, P., Asnicar, F., Dubois, L., Golzato, D., Armanini,
1375 F., Cumbo, F., Huang, K.D., Manara, S., et al. (2023). The person-to-person transmission landscape of
1376 the gut and oral microbiomes. *Nature* 614, 125–135. <https://doi.org/10.1038/s41586-022-05620-1>.
- 1377 95. Liu, Q., Du, X., Hong, X., Li, T., Zheng, B., He, L., Wang, Y., Otto, M., and Li, M. (2015). Targeting
1378 Surface Protein SasX by Active and Passive Vaccination To Reduce *Staphylococcus aureus*
1379 Colonization and Infection. *Infect Immun* 83, 2168–2174. <https://doi.org/10.1128/iai.02951-14>.
- 1380 96. Huttenhower, C., Gevers, D., Knight, R., Abubucker, S., Badger, J.H., Chinwalla, A.T., Creasy,
1381 H.H., Earl, A.M., FitzGerald, M.G., Fulton, R.S., et al. (2012). Structure, function and diversity of the
1382 healthy human microbiome. *Nature* 486, 207–214. <https://doi.org/10.1038/nature11234>.
- 1383 97. Bäckhed, F., Roswall, J., Peng, Y., Feng, Q., Jia, H., Kovatcheva-Datchary, P., Li, Y., Xia, Y., Xie,
1384 H., Zhong, H., et al. (2015). Dynamics and Stabilization of the Human Gut Microbiome during the First
1385 Year of Life. *Cell Host Microbe* 17, 690–703. <https://doi.org/10.1016/j.chom.2015.04.004>.
- 1386 98. Hyatt, D., Chen, G.-L., LoCascio, P.F., Land, M.L., Larimer, F.W., and Hauser, L.J. (2010).
1387 Prodigal: prokaryotic gene recognition and translation initiation site identification. *Bmc Bioinformatics*
1388 11, 119. <https://doi.org/10.1186/1471-2105-11-119>.

- 1389 99. Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K., and Madden, T.L.
1390 (2009). BLAST+: architecture and applications. *Bmc Bioinformatics* 10, 421.
1391 <https://doi.org/10.1186/1471-2105-10-421>.
- 1392 100. Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2 – Approximately Maximum-Likelihood
1393 Trees for Large Alignments. *Plos One* 5, e9490. <https://doi.org/10.1371/journal.pone.0009490>.
- 1394 101. Enright, A.J., Dongen, S.V., and Ouzounis, C.A. (2002). An efficient algorithm for large-scale
1395 detection of protein families. *Nucleic Acids Res* 30, 1575–1584. <https://doi.org/10.1093/nar/30.7.1575>.
- 1396 102. García-Bayona, L., and Comstock, L.E. (2019). Streamlined Genetic Manipulation of Diverse
1397 *Bacteroides* and *Parabacteroides* Isolates from the Human Gut Microbiota. *Mbio* 10, e01762-19.
1398 <https://doi.org/10.1128/mbio.01762-19>.
- 1399 103. Ito, T., Gallegos, R., Matano, L.M., Butler, N.L., Hantman, N., Kaili, M., Coyne, M.J., Comstock,
1400 L.E., Malamy, M.H., and Barquera, B. (2020). Genetic and Biochemical Analysis of Anaerobic
1401 Respiration in *Bacteroides fragilis* and Its Importance In Vivo. *Mbio* 11, e03238-19.
1402 <https://doi.org/10.1128/mbio.03238-19>.
- 1403 104. Lim, B., Zimmermann, M., Barry, N.A., and Goodman, A.L. (2017). Engineered Regulatory
1404 Systems Modulate Gene Expression of Human Commensals in the Gut. *Cell* 169, 547-558.e15.
1405 <https://doi.org/10.1016/j.cell.2017.03.045>.
- 1406 105. Ahuja, U., Shokeen, B., Cheng, N., Cho, Y., Blum, C., Coppola, G., and Miller, J.F. (2016).
1407 Differential regulation of type III secretion and virulence genes in *Bordetella pertussis* and *Bordetella*
1408 *bronchiseptica* by a secreted anti- σ factor. *Proc National Acad Sci* 113, 2341–2348.
1409 <https://doi.org/10.1073/pnas.1600320113>.
- 1410 106. Thein, M., Sauer, G., Paramasivam, N., Grin, I., and Linke, D. (2010). Efficient Subfractionation of
1411 Gram-Negative Bacteria for Proteomics Studies. *J Proteome Res* 9, 6135–6147.
1412 <https://doi.org/10.1021/pr1002438>.
- 1413 107. Hughes, C.S., Moggridge, S., Müller, T., Sorensen, P.H., Morin, G.B., and Krijgsveld, J. (2019).
1414 Single-pot, solid-phase-enhanced sample preparation for proteomics experiments. *Nat. Protoc.* 14, 68–
1415 85. <https://doi.org/10.1038/s41596-018-0082-x>.

1416 108. Cox, J., and Mann, M. (2008). MaxQuant enables high peptide identification rates, individualized
1417 p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat Biotechnol* 26, 1367–1372.
1418 <https://doi.org/10.1038/nbt.1511>.

1419 109. Jimenez-Morales, D., Campos, A.R., Dollen, J.V., Krogan, N., and Swaney, D. artMS: Analytical R
1420 tools for Mass Spectrometry. <https://bioconductor.org/packages/release/bioc/html/artMS.html>.

1421 110. Goedhart, J., and Luijsterburg, M.S. (2020). VolcaNoseR is a web app for creating, exploring,
1422 labeling and sharing volcano plots. *Sci. Rep.* 10, 20560. <https://doi.org/10.1038/s41598-020-76603-3>.

1423 111. Fields, B., Moeskjær, S., Friman, V., Andersen, S.U., and Young, J.P.W. (2021). MAUI-seq:
1424 Metabarcoding using amplicons with unique molecular identifiers to improve error correction. *Mol Ecol*
1425 *Resour* 21, 703–720. <https://doi.org/10.1111/1755-0998.13294>.

1426 112. Bushnell, B., Rood, J., and Singer, E. (2017). BBMerge – Accurate paired shotgun read merging
1427 via overlap. *Plos One* 12, e0185056. <https://doi.org/10.1371/journal.pone.0185056>.

1428 113. Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., Beghini, F., Manghi, P.,
1429 Tett, A., Ghensi, P., et al. (2019). Extensive Unexplored Human Microbiome Diversity Revealed by
1430 Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* 176, 649-
1431 662.e20. <https://doi.org/10.1016/j.cell.2019.01.001>.

1432 114. Steinegger, M., and Söding, J. (2017). MMseqs2 enables sensitive protein sequence searching for
1433 the analysis of massive data sets. *Nat. Biotechnol.* 35, 1026–1028. <https://doi.org/10.1038/nbt.3988>.

1434 115. Paul, B.G., Burstein, D., Castelle, C.J., Handa, S., Arambula, D., Czornyj, E., Thomas, B.C.,
1435 Ghosh, P., Miller, J.F., Banfield, J.F., et al. (2017). Retroelement guided protein diversification abounds
1436 in vast lineages of bacteria and archaea. *Nat Microbiol* 2, 17045–17045.
1437 <https://doi.org/10.1038/nmicrobiol.2017.45>.

1438 116. Langmead, B., and Salzberg, S.L. (2012). Fast gapped-read alignment with Bowtie 2. *Nat*
1439 *Methods* 9, 357–359. <https://doi.org/10.1038/nmeth.1923>.

1440 117. Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler
1441 transform. *Bioinformatics* 25, 1754–1760. <https://doi.org/10.1093/bioinformatics/btp324>.

1442 118. Edgar, R.C. (2004). MUSCLE: a multiple sequence alignment method with reduced time and
1443 space complexity. *Bmc Bioinformatics* 5, 113. <https://doi.org/10.1186/1471-2105-5-113>.

1444 119. Hall, M., Hasegawa, Y., Yoshimura, F., and Persson, K. (2018). Structural and functional
1445 characterization of shaft, anchor, and tip proteins of the Mfa1 fimbria from the periodontal pathogen
1446 *Porphyromonas gingivalis*. *Sci. Rep.* 8, 1793. <https://doi.org/10.1038/s41598-018-20067-z>.

1447 120. Shimoyama, Y. pyGenomeViz: A genome visualization python package for comparative genomics.
1448 <https://github.com/moshi4/pyGenomeViz>.

1449 121. Chen, S. (2023). Ultrafast one-pass FASTQ data preprocessing, quality control, and deduplication
1450 using fastp. *iMeta* 2, e107. <https://doi.org/10.1002/imt2.107>.

1451 122. Md, V., Misra, S., Li, H., and Aluru, S. (2019). Efficient Architecture-Aware Acceleration of BWA-
1452 MEM for Multicore Systems. 2019 IEEE Int. Parallel Distrib. Process. Symp. (IPDPS) 00, 314–324.
1453 <https://doi.org/10.1109/ipdps.2019.00041>.

1454 123. Barnett, D.W., Garrison, E.K., Quinlan, A.R., Strömberg, M.P., and Marth, G.T. (2011). BamTools:
1455 a C++ API and toolkit for analyzing and managing BAM files. *Bioinformatics* 27, 1691–1692.
1456 <https://doi.org/10.1093/bioinformatics/btr174>.

1457 124. Ewels, P., Magnusson, M., Lundin, S., and Käller, M. (2016). MultiQC: summarize analysis results
1458 for multiple tools and samples in a single report. *Bioinformatics* 32, 3047–3048.
1459 <https://doi.org/10.1093/bioinformatics/btw354>.

1460 125. Liao, Y., Smyth, G.K., and Shi, W. (2014). featureCounts: an efficient general purpose program for
1461 assigning sequence reads to genomic features. *Bioinformatics* 30, 923–930.
1462 <https://doi.org/10.1093/bioinformatics/btt656>.

1463 126. Tarazona, S., Furió-Tarí, P., Turrà, D., Pietro, A.D., Nueda, M.J., Ferrer, A., and Conesa, A.
1464 (2015). Data quality aware analysis of differential expression in RNA-seq with NOISeq R/Bioc package.
1465 *Nucleic Acids Res.* 43, e140–e140. <https://doi.org/10.1093/nar/gkv711>.

1466