

1

2 **The Landscape of Type VI Secretion across Human Gut Microbiomes**
3 **Reveals its Role in Community Composition**

4

5 Adrian J. Verster^{1,*}, Benjamin D. Ross^{2,*}, Matthew C. Radey², Yiqiao Bao^{3,4}, Andrew L.
6 Goodman^{3,4}, Joseph D. Mougous^{2,5,†}, Elhanan Borenstein^{1,6,7,†}

7

8

9

10 ¹Department of Genome Sciences, University of Washington,
11 Seattle, WA, 98195, USA

12

13 ²Department of Microbiology, School of Medicine, University of Washington,
14 Seattle, WA 98195, USA

15

16 ³Department of Microbial Pathogenesis, Yale University School of Medicine,
17 New Haven, CT 06510, USA

18

19 ⁴Microbial Sciences Institute, Yale University School of Medicine,
20 West Haven, CT 06516, USA

21

22 ⁵Howard Hughes Medical Institute, School of Medicine, University of Washington,
23 Seattle, WA 98195, USA

24

25 ⁶Department of Computational Science and Engineering, University of Washington,
26 Seattle, WA 98195, USA

27

28 ⁷Santa Fe Institute, Santa Fe, NM 87501 USA

29

30 *Equal contribution

31

32 †To whom correspondence should be addressed: elbo@uw.edu, mougous@uw.edu

33

1 **Abstract**

2 While the composition of the human gut microbiome has been well defined, the forces
3 governing its assembly are poorly understood. Recently, prominent members of this
4 community from the order *Bacteroidales* were shown to possess the type VI secretion
5 system (T6SS), which mediates contact-dependent antagonism between Gram-negative
6 bacteria. However, the distribution of the T6SS in human gut microbiomes and its role
7 have not yet been characterized. To address this challenge, we construct an extensive
8 catalog of T6SS effector/immunity (E-I) genes from three genetic architectures (GA1-3)
9 found in *Bacteroidales* genomes. We then use metagenomic analysis to assess the
10 abundances of these genes across a large set of gut microbiome samples. We find that
11 despite E-I diversity across reference strains, each individual microbiome harbors a
12 limited set of E-I genes representing a single E-I genotype. Importantly, for GA1-2, these
13 genotypes are not associated with a specific species, suggesting selection for
14 compatibility. GA3, in contrast, is restricted to *B. fragilis*, and its low diversity reflects a
15 single *B. fragilis* strain per sample. We further show that in infant microbiomes GA3 is
16 enriched and *B. fragilis* strains are replaced over time, suggesting competition for
17 dominance in developing microbiomes. Finally, we find a strong association between the
18 presence of GA3 and increased abundance of *Bacteroides*, indicating that this system
19 confers a selective advantage *in vivo* in *Bacteroides* rich ecosystems. Combined, our
20 findings provide the first comprehensive characterization of the T6SS landscape in the
21 human microbiome, implicating it in both intra- and inter-species interactions.

22

1 **Introduction**

2 Bacterial communities are of fundamental importance to natural ecosystems (Prosser et
3 al. 2007). While cooperative interactions between the species comprising such
4 communities are common (Morris et al. 2012), it is clear that bacteria in these settings
5 also experience pervasive antagonism from surrounding cells (Hibbing et al. 2010).
6 Indeed, the genomes of bacteria encode a wealth of dedicated interbacterial antagonism
7 pathways (Zhang et al. 2012). Some of these function through the production of diffusible
8 small molecules (Riley and Wertz 2002), whereas others utilize proteinaceous toxins. A
9 prevalent pathway mediating the transfer of toxic proteins between bacteria is the type VI
10 secretion system (T6SS). This system has been most thoroughly studied in
11 *Proteobacteria*, though it is found in several phyla of Gram-negative bacteria (Russell et
12 al. 2014a).

13 The T6S apparatus transfers toxic effector proteins from donor to recipient cells
14 by a mechanism dependent upon cell contact (Russell et al. 2014a). Characterized
15 effector proteins are thus far without exception enzymes that target conserved, essential
16 features of the bacterial cell, such as peptidoglycan, phospholipids, and nucleic acids.
17 This feature of effector proteins, taken together with the fact that T6SS targeting is not
18 dependent on a specific receptor, confers broad activity against Gram-negative cells.
19 Indiscriminate effector transfer also extends to kin cells; therefore, cells with the T6SS
20 produce immunity proteins that inactivate cognate toxins through active site occlusion
21 (Whitney et al. 2015).

22 Given its wide phylogenetic distribution and its capacity to target diverse recipient
23 cells, the T6SS is likely to play an important role in the assembly and composition of

1 bacterial communities. Indeed, there are recent reports consistent with the pathway
2 mediating bacterial interaction in environmental communities. For instance, T6S genes
3 were found to be enriched and under positive selection in the barley rhizosphere
4 (Bulgarelli et al. 2015), and T6S phospholipase effectors were detected in metagenomes
5 from diverse sources (Egan et al. 2015). To date, however, systematic studies of the
6 impact of T6SS on microbial community assembly are lacking.

7 The human gut microbiome is a dense ecosystem whose composition is
8 paramount to its function (Walter and Ley 2011). Factors such as diet, immune status,
9 and host genetics have each been implicated in shaping the gut community (Kau et al.
10 2011), yet the contribution of direct interbacterial competition to the structure of this
11 community remains poorly understood. Recently, a T6SS-like pathway was detected in
12 *Bacteroidetes*, the most abundant Gram-negative phylum in the human gut (Coyne et al.
13 2014; Russell et al. 2014b). Additional work demonstrated that T6S contributes to the
14 fitness of *Bacteroides fragilis* in competition with other bacteria *in vitro* and in
15 gnotobiotic mice (Russell et al. 2014b; Chatzidaki-Livanis et al. 2016; Hecht et al. 2016;
16 Wexler et al. 2016). These and other data show that the mammalian GI tract is physically
17 conducive to T6S-dependent interbacterial antagonism, suggesting a potential impact of
18 this pathway on the composition of the human gut microbiome (Dong et al. 2013; Sana et
19 al. 2016). Here, we sought to define the distribution of *Bacteroidales* T6S and to explore
20 its function in the human gut microbiome through the analysis of several publicly
21 available metagenomic datasets. These datasets allow us to study the outcome of natural
22 community dynamics in the gut microbiome, and we reasoned that their analysis could
23 therefore provide unique insight into the physiologic role of T6SS-dependent competition

1 in this ecosystem. Our findings reveal the prevalence of this pathway in intact human gut
2 microbial communities, highlight striking and non-random patterns in its distribution
3 across samples, and suggest an active role for T6S in intra- and inter-species interactions
4 in the gut.

5 **Results**

6 **Detection of T6S E–I pairs in the human gut microbiome**

7 We first set out to characterize the prevalence and distribution of T6SS genes in the gut
8 microbiomes of healthy adult individuals. Based on their organization and content,
9 *Bacteroidales* T6S gene clusters can be divided into three distinct subtypes, termed
10 genetic architecture 1-3 (GA1-3) (Coyne et al. 2016). Each T6S subtype possesses one or
11 more cassettes at stereotyped positions that contain variable genes predicted or
12 demonstrated to encode effector–immunity (E–I) pairs (Russell et al. 2014b; Chatzidaki-
13 Livanis et al. 2016; Coyne et al. 2016; Wexler et al. 2016). As T6S-based antagonism is
14 determined by the effector and immunity genes of donor and recipient cells, respectively,
15 the identification of E–I pairs provides information regarding the potential for
16 interbacterial interactions mediated by this system (Hood et al. 2010). Furthermore, since
17 these cassettes are variable within, but unique among the T6S subtypes, estimation of the
18 abundance of these genes within metagenomes can serve as a proxy for the presence and
19 distribution of GA1-3.

20 To define the E–I repertoire associated with GA1-3, we searched genes within T6-
21 associated variable cassettes from *Bacteroidales* genomes for those with hallmarks of
22 known T6 effector and immunity factors. These included fusion to modular adaptor
23 domains, reduced GC content, bicistronic arrangement, and similarity to protein families

1 defined by their association with characterized E–I pairs (see Methods for a complete
 2 description of annotation criteria; Figures 1A and S1). In total, we identified 9 GA1, 18
 3 GA2, and 14 GA3 putative E–I pairs. As expected, GA3 pairs were identified only in *B.*
 4 *fragilis*, whereas GA1 and GA2 pairs were detected throughout the order. Importantly,
 5 we did not identify homologs of GA1-3 E–I genes outside of *Bacteroidales*.
 6

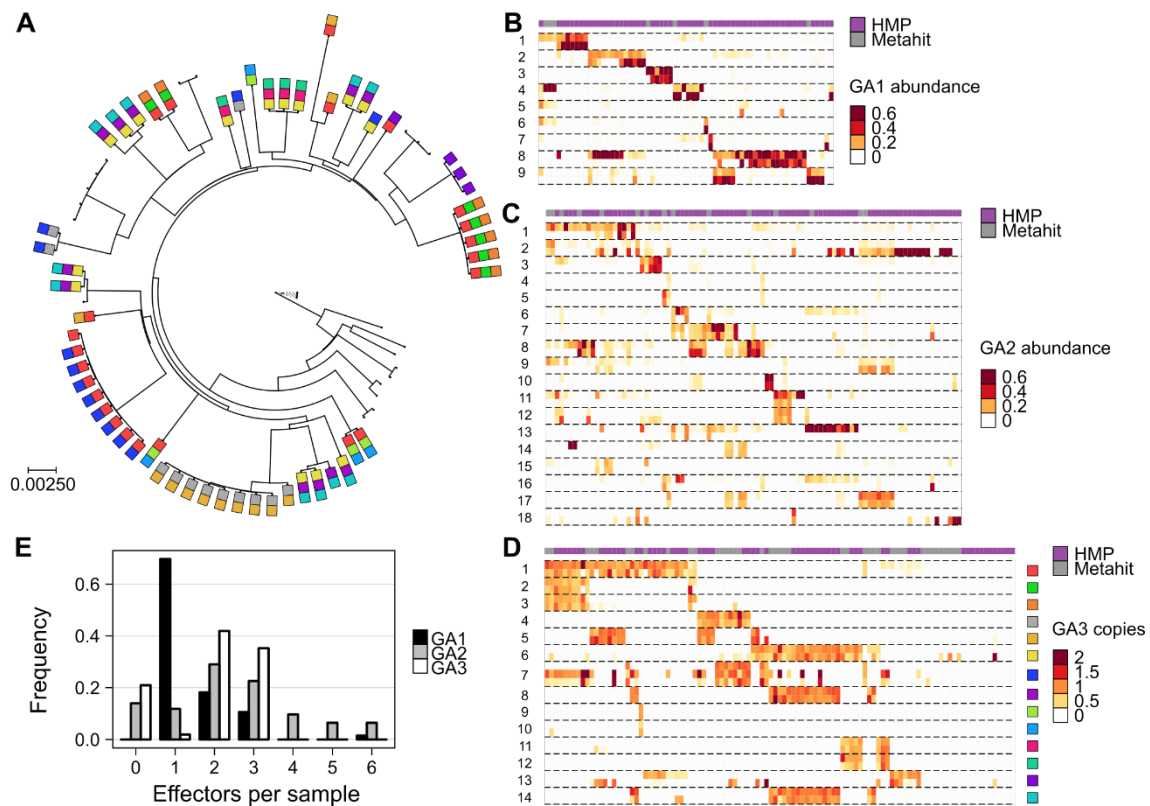


Figure 1 | *Bacteroides* T6SS E–I genes are abundant in human gut microbiome samples.

A. A maximum likelihood phylogeny of *B. fragilis* reference strains constructed from concatenated marker genes. Phylogenetic distance is measured as substitutions per site on the marker genes. GA3 effector genes are represented as colored squares (using the same color coding as in panel D). **B–D.** Each heatmap illustrates the abundance of E–I genes for one of the T6SS subsystems. Each row corresponds to a different E–I pair (effector, top; immunity, bottom). Columns represent the samples analyzed (HMP, purple; MetaHIT, gray). For GA1 and GA2, only samples in which at least 100 reads mapped to the E–I genes of a given subsystem are included, and abundance is measured as fraction of the total abundance of E–I genes in a given sample. For GA3, only samples in which *B. fragilis* is present are included and E–I abundance is normalized by the abundance of *B. fragilis*-specific marker genes, hence measuring the average number of copies per *B. fragilis* genome. **E.** Histograms showing the number of effector genes detected (at >10% of the most abundant effector gene) in each sample.

1 To estimate the abundance of these E–I pairs in gut microbiomes, we obtained
2 metagenomic datasets derived from healthy donor samples of the HMP (Human
3 Microbiome Project 2012) and MetaHIT (Qin et al. 2010) studies. We next mapped the
4 reads from each sample to the E–I genes using a sequence identity threshold
5 demonstrated to maintain E–I compatibility (Unterweger et al. 2014). Our results indicate
6 that T6S is prevalent in the human gut microbiome; of the 246 samples analyzed, we
7 detected E–I genes in 166 (67%). Moreover, each E–I pair in our list was detected in at
8 least one microbiome sample, with an average of 12.2 occurrences.

9

10 **T6S E–I pairs display low diversity within human gut microbiome samples**

11 The systematic characterization of E–I abundance in metagenomic samples provided a
12 unique opportunity to examine the distribution of the genes associated with each genetic
13 architecture across healthy gut microbiomes. We first focused on GA1 and GA2, which
14 utilize unique complements of effectors, but share the ability to undergo conjugative
15 transfer between species belonging to the order *Bacteroidales* (Coyne et al. 2016).
16 Surprisingly, we found that the complement of GA1- and GA2-associated E–I genes in a
17 typical microbiome is small, with only few pairs per sample, comparable to the number
18 of pairs usually detected in a single genome (Figure 1B,C, E). Moreover, in many cases,
19 the same complement of E–I genes was detected in multiple samples. Henceforth, we
20 refer to these combinations as E–I genotypes. This pattern suggests that either each
21 sample is dominated by a single strain that harbors the observed E–I genotype or that
22 there exists selective pressure for compatible E–I genes across multiple strains or species
23 in a sample.

1 To further explore these possibilities, we focused our attention on the most
2 prominent members of the genus *Bacteroides*. Other genera in the order *Bacteroidales* are
3 less abundant constituents of the microbiome and based on reference genomes do not
4 often harbor GA1 or GA2. We identified a set of species-specific single-copy marker
5 genes for each *Bacteroides* species and estimated their abundance in each sample. Next
6 we compared marker gene abundance to that of GA1 and GA2 E–I genes across samples
7 (see Methods). We found that the abundance of these E–I genes was not consistent with
8 that of any individual species (Figures 2A and S2). These findings suggest that multiple
9 species co-existing in a microbiome typically encode a single GA1 and/or GA2 E–I
10 genotype, potentially due to selective pressure for maintenance of E–I compatibility.

11 We next examined GA3 E–I genes, and found, as in GA1 and GA2, that each
12 sample harbors only a small set of E–I pairs (Figure 1D,E). Moreover, observed GA3 E–I
13 genotypes matched those detected in reference genomes (Figure S3), and appeared
14 randomly distributed between the American (HMP) and European (MetaHIT) datasets
15 (Figure 1D). However, in contrast to GA1 and GA2, we found a strong correlation ($R =$
16 0.94) between the abundance of GA3 effector and immunity genes and that of a single
17 species, *B. fragilis* (Figures 2A-B and S2). This finding is supported by our curation of
18 E–I genes and by previous studies (Coyne et al. 2016) that found GA3 restricted to *B.*
19 *fragilis*. Importantly, however, our finding confirms that this restriction of GA3 to *B.*
20 *fragilis* observed in sequenced reference genomes holds across naturally occurring
21 communities.

1 We hypothesized that the pattern of GA3 E-I genotypes we observed could be
2 explained by the dominance of a single *B. fragilis* strain within each individual
3 microbiome. Indeed, prior studies suggest that *B. fragilis* exhibits relatively low diversity
4 within individuals (Yassour et al. 2016). To confirm that this pattern is also observed in

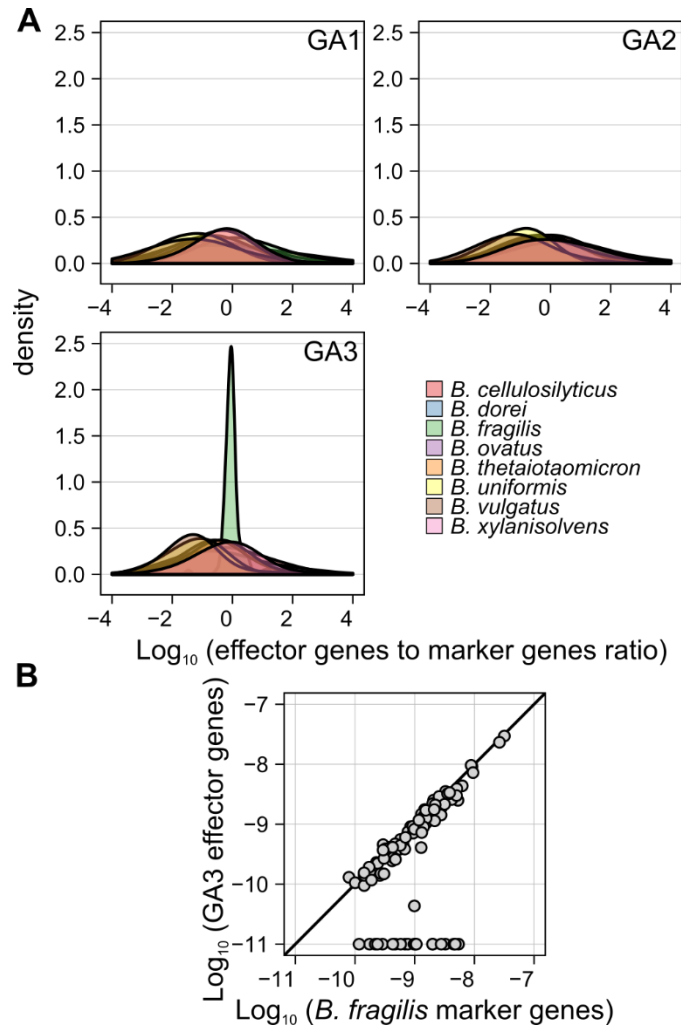


Figure 2 | Differential associations between T6SS and *Bacteroides* spp.

A. Density plots showing the distribution across samples of the ratio between the average abundance of detected effector genes from each T6SS subsystem and the average abundance of species-specific marker genes for different *Bacteroides* spp. Only samples in which at least 100 reads mapped to the E-I genes of a given subsystem and only species for which at least 5 genomes were available (and therefore marker genes can be robustly inferred) are included. **B.** Scatter plot of the average abundance of detected GA3 effector genes vs. the average abundance of *B. fragilis*-specific marker genes. Only samples in which *B. fragilis* is present are included. A small factor (10^{-11}) has been added to each abundance value to allow transformation to a logarithmic scale.

1 HMP and MetaHIT samples, we first measured nucleotide diversity in species-specific
2 markers of *Bacteroides* spp. We found that, within an individual, *B. fragilis* possesses the
3 lowest average SNP diversity of well represented members of the *Bacteroides* genus
4 (Figure S4A). We then used a previously developed method for inferring the most likely
5 set of strains in metagenomic samples based on nucleotide variants, combined with a
6 phylogenetic analysis of these inferred strains to determine the number of different
7 monophyletic groups of strains present in each sample (see Methods). We found that *B.*
8 *fragilis* inferred strains in every HMP and MetaHIT sample formed a single
9 monophyletic group, indicating that extant *B. fragilis* strains in each sample are likely
10 derived from a single colonization event. Moreover, the set of E–I pairs detected in each
11 metagenomic sample generally matched the set of E–I pairs found in the reference strains
12 closest in the phylogenetic tree to the inferred strains, especially when the distance of
13 inferred strains to their nearest reference was low (Figure S5).

14 To experimentally confirm our computational findings, we additionally selected
15 20 randomized *B. fragilis* colonies isolated from two healthy adults and subjected these to
16 whole genome sequencing. Consistent with our findings using metagenomic data, our
17 sequencing showed that a single clonal strain of *B. fragilis* dominates the microbiome of
18 these individuals (Figure S4B).

19

20 ***B. fragilis* GA3 is important in the developing microbiome**

21 The finding that the presence of singular GA3 genotypes within individuals is due to the
22 dominance of one *B. fragilis* strain motivated us to investigate the role of this system in
23 the microbiome. We reasoned that in this dense and competitive microbiome ecosystem

1 (Ley et al. 2006), an antagonistic pathway such as the T6SS might provide a fitness
2 advantage. The system could mediate antagonism against other *B. fragilis* strains closely
3 related organisms such as other *Bacteroides* spp., other Gram-negative inhabitants of the
4 microbiome, or a combination of these (Chatzidaki-Livanis et al. 2016; Hecht et al. 2016;
5 Wexler et al. 2016). Assuming such a role for T6SS, we further reasoned that in the
6 microbiome of infants, which is less stable than that of adults(Sharon et al. 2013), the
7 function of an antagonistic pathway like the T6SS might be more pronounced. To test this
8 hypothesis, we obtained publically available metagenomic datasets derived from infant
9 gut microbiomes (Backhed et al. 2015; Kostic et al. 2015; Vatanen et al. 2016; Yassour et
10 al. 2016). We then identified samples that contain *B. fragilis* but lack GA3-associated
11 structural genes (see Methods) in both adult and infant datasets. Such samples indicate
12 the presence of *B. fragilis* strains unable to intoxicate competitor bacteria using this
13 pathway. The availability of well-assembled *B. fragilis* reference genomes that lack the
14 GA3 gene cluster lent credence to this approach (Wexler et al. 2016). We found that
15 infant microbiomes containing *B. fragilis* are significantly less likely to lack GA3-
16 associated structural genes relative to those of adults (Fisher's exact test, $P < 0.01$, 8%
17 infants, 23% adults; $n = 276$; Figures 3A and S6).

18 This finding suggests that GA3 provides an advantage for *B. fragilis* in early life;
19 however, the selection pressure underlying this advantage remained unclear. Several
20 independent studies using gnotobiotic mice have shown that the GA3 T6SS can play a
21 major role in the competition between *B. fragilis* strains in the gut (Chatzidaki-Livanis et
22 al. 2016; Hecht et al. 2016; Wexler et al. 2016). However, *B. fragilis* is thought to be
23 stable after acquisition from the mother, and inter-strain competition within the human

1 gut microbiome has not been documented for this organism (Faith et al. 2013; Nayfach et
 2 al. 2016). Aiming to capture such processes in the developing microbiome, we estimated
 3 the abundance of GA3 E-I genes for individual infant samples as we did for adults. In
 4 general, the E-I landscape of infants mirrors that of adults, with generally a single
 5 genotype present in each sample (Figure S7). Moreover, many of the most prevalent E-I
 6 genotypes we observed in adults are also frequent in infants.

7 Notably, the infant microbiome datasets we analyzed include multiple samples per

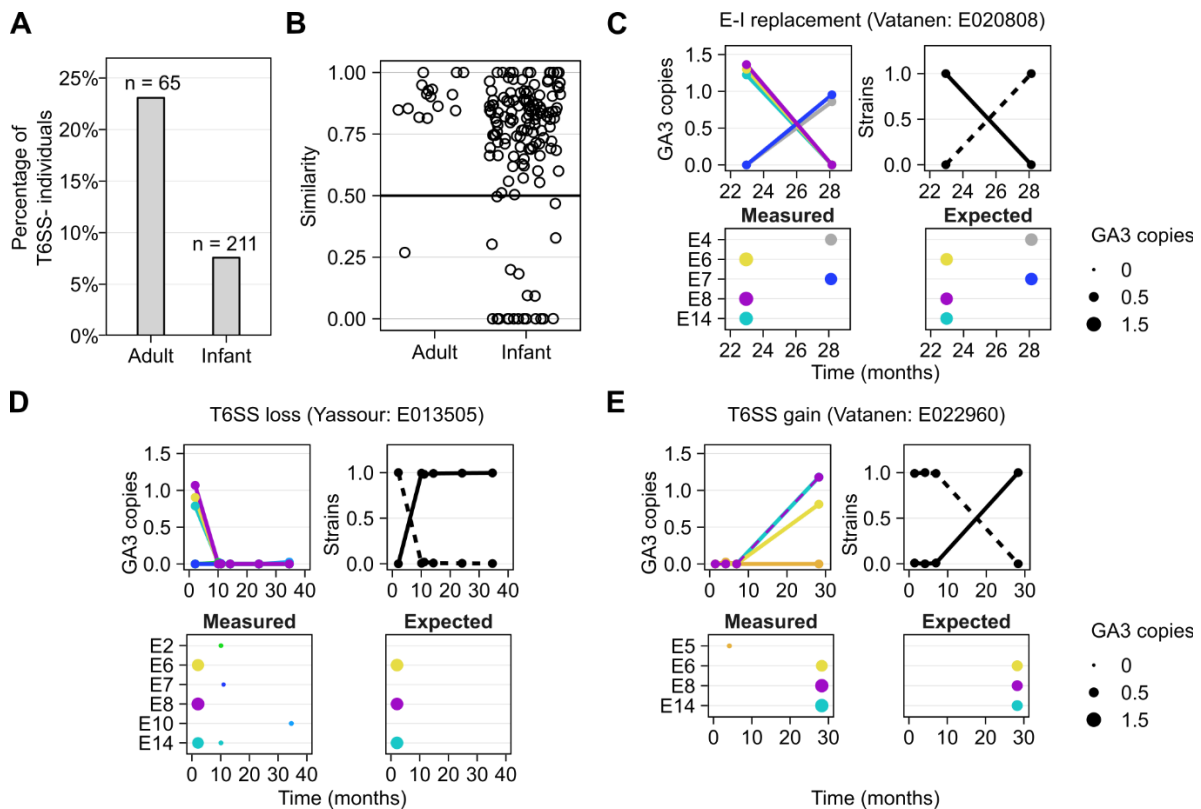


Figure 3 | E-I turnover and strain replacement in infant microbiomes.

A. The percentage of individuals of those harboring *B. fragilis*, that lack the GA3 T6SS across adult and infant datasets. **B.** The minimal similarity (measured by the Jaccard similarity coefficient) in GA3 E-I gene content between the first time point and every subsequent time point in adults and infants. **C-E.** Examples of E-I turnover events and corresponding strain replacement events are shown. The plots on the upper and bottom left in each panel illustrate the estimated abundance of GA3 effector genes (measured as copies per *B. fragilis* genome) over time, with the plot on the upper right illustrating the estimated frequency of inferred strains in these samples. Only samples in which *B. fragilis* is present are shown. The bottom right plot illustrates the expected abundance of the various effector genes based on the effector genes encoded by reference strains that are phylogenetically close to the inferred strains.

1 individual, thereby allowing us to examine the temporal dynamics of *B. fragilis* and of
2 T6SS genes. Surprisingly, this analysis revealed many instances in which the E-I
3 genotype of an individual has changed between samples (Figure 3B). In total, we
4 observed E-I turnover in 22 of the 117 infants for which longitudinal data was available.
5 Such E-I turnover events include instances where one GA3 genotype is replaced by
6 another (Figure 3C), but also gains and losses of the T6SS (Figure 3D-E). To further
7 confirm these E-I dynamics, we used the strain inference method described above. We
8 detected a corresponding strain replacement in 17 of the 22 individuals in which an E-I
9 turnover event was observed (Figures 3C-E and S8). Moreover, comparing the set of E-I
10 genes detected in each sample to those encoded by the reference strains phylogenetically
11 closest to the inferred strain, we further find overall agreement between observed and
12 expected E-I turnover events. Notably, instances of one strain replaced by another with a
13 similar E-I genotype (Figure S8; Vatanen:T014827) or of transient co-existence of E-I
14 genotypes (Figure S8; Backhed:587) were also observed. Interestingly, examining the
15 few HMP adult individuals for which data was available from multiple visits, we found
16 one adult in which the E-I genotype similarly changed over time (Figure 3B).

17

18 ***B. fragilis* GA3 T6SS is associated with shifts in community composition**

19 Due to its lower frequency in adult microbiomes compared to those of infants, the GA3
20 T6SS is absent in many adult samples in which *B. fragilis* can be detected (23%), offering
21 a unique opportunity to compare the community composition in samples with or without
22 GA3, and to uncover community-associated outcomes of T6S activity in the human gut.

1 To this end, we obtained the taxonomic profile of all HMP samples (see Methods) and
2 identified associations between these profiles and the presence of T6SS structural genes.
3 We first compared overall community composition between samples as measured by the
4 Bray-Curtis distance. We found that samples harboring *B. fragilis* and GA3 genes
5 (T6SS+) significantly differ in community composition from samples harboring *B.*
6 *fragilis* but lacking these genes (T6SS-; $P < 0.01$ PERMANOVA; $n = 51$). Examining the
7 abundance of each genera across samples, we further identified four genera whose
8 abundance in T6SS+ versus T6SS- samples significantly differs (Wilcoxon rank sum test;
9 $FDR < 0.05$; Figure 4 and Table S1). Specifically, we found that the abundance of
10 *Bacteroides* is positively correlated with the presence of GA3, which is consistent with
11 experimental and theoretical work indicating that members of this genus are most likely

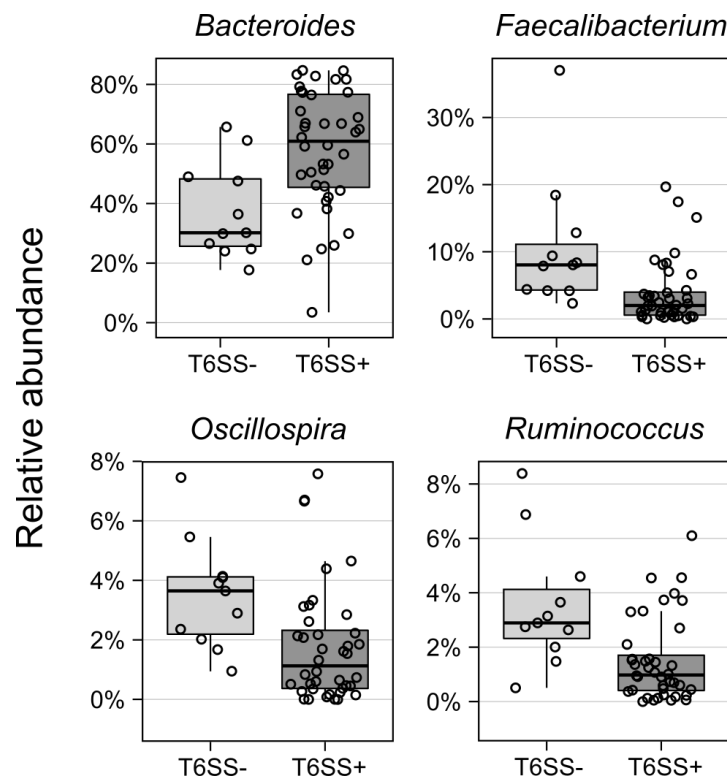


Figure 4 | Differentially abundant genera between T6SS+ and T6SS- HMP samples. Abundances are based on a 16S rRNA survey and only genera whose abundances are significantly different in T6SS+ vs. T6SS- samples (at $FDR < 0.05$) are plotted.

1 to compete with *B. fragilis* for its niche (Russell et al. 2014b; Trosvik and de Muinck
2 2015; Chatzidaki-Livanis et al. 2016; Hecht et al. 2016; Wexler et al. 2016). Furthermore,
3 the genera *Faecalibacterium*, *Oscillospira* and *Ruminococcus* from the phylum
4 *Firmicutes* were negatively correlated with GA3. Gram-positive organisms are not targets
5 of any known T6SS; therefore, the observed decreases in abundance of these genera in
6 T6SS+ microbiomes are likely to be the indirect result of selection for GA3 occurring in
7 communities with an increased ratio of *Bacteroidetes* to *Firmicutes*.

8

9 **Discussion**

10 Despite the wide distribution of T6S in Gram-negative bacteria, little is known about its
11 role in natural communities. Here, we surveyed human microbiome samples and
12 discovered that these communities are replete with bacteria containing the T6SS. We
13 focused our analyses on genes specific to the order *Bacteroidales*; thus, our findings are
14 an underestimate of the prevalence and impact of this system in gut communities.
15 Nonetheless, our characterization of T6SS E–I gene distribution in human gut
16 microbiomes suggests that this contact-dependent pathway plays a role in competition
17 and selection at multiple levels.

18 We observed a markedly low diversity of T6SS E–I genes in human microbiome
19 samples. Specifically, a single genotype of GA1 and GA2 E–I genes is found in each
20 microbiome, yet the abundance of these genes does not correlate with that of any one
21 species in the *Bacteroides* genus. This supports a model in which antagonism via GA1
22 and GA2 exerts selective pressure for compatibility between *Bacteroides* spp. in the gut.
23 We postulate that in the case of GA1 and GA2, horizontal transfer facilitates E–I

1 compatibility. Indeed, Comstock and colleagues found that transfer of GA1 and GA2 can
2 occur between species within the microbiome of an individual (Coyne et al. 2014).

3 Similar to GA1 and GA2, we found a single genotype of GA3 in each
4 microbiome, but were able to demonstrate that GA3 is restricted to *B. fragilis* and can be
5 explained by the presence of a single *B. fragilis* strain in each individual. This is
6 consistent with experimental studies demonstrating a role for GA3-dependent
7 competition between *B. fragilis* strains (Chatzidaki-Livanis et al. 2016; Hecht et al. 2016;
8 Wexler et al. 2016). We found, however, that infant microbiomes commonly exhibit a
9 strain replacement accompanied by E–I turnover.

10 While the low number of adult longitudinal samples prohibits us from statistically
11 comparing strain replacement rates in infants and adults, the observed strain replacements
12 in infants and our finding that *B. fragilis* lacking the GA3 T6SS is more common in
13 adults suggest that, early in life, *B. fragilis* strains compete for dominance, potentially
14 contributing to the establishment of this system in the infant microbiome. This is further
15 supported by the fact that we rarely find co-existing strains of *B. fragilis*, and that typical
16 turnover events are characterized by the rapid transition of one dominant strain to
17 another. Indeed, gnotobiotic experiments indicate that inter-strain competition is strong in
18 mice colonized solely with *B. fragilis*, and diminished when the relative abundance of *B.*
19 *fragilis* is reduced in a more complex community (Wexler et al. 2016). Therefore, the
20 infant microbiome may represent a particularly dynamic ecosystem in which the GA3
21 T6SS facilitates *B. fragilis* strain competition.

22 We find that strains of *B. fragilis* lacking GA3 are more commonly found in
23 adults than infants. This could arise either by the replacement of T6SS+ with T6SS-

1 strains, or by the loss of the T6SS system from a previously T6SS+ strain of *B. fragilis*.
2 This decline in *B. fragilis* GA3 prevalence in adulthood may reflect a change in its
3 selective advantage. Indeed, there is precedent for the lability of T6S in bacteria
4 undergoing strong shifts in environmental context, such as *Burkholderia mallei* and
5 *Bordetella* spp. (Schwarz et al. 2010). It is likely that community effects buffering the *B.*
6 *fragilis* niche develop with the maturation of the more stable adult gut community.
7 Stabilization over development appears to render GA3 dispensable in certain contexts, for
8 instance within those microbiomes that contain lower populations of potential *B. fragilis*
9 competitors and known targets of the GA3 pathway, other *Bacteroides* spp. Our findings
10 and others suggest that *B. fragilis* primarily faces selective pressure from closely-related
11 species (Chatzidaki-Livanis et al. 2014; Russell et al. 2014b; Chatzidaki-Livanis et al.
12 2016; Hecht et al. 2016; Roelofs et al. 2016; Wexler et al. 2016). Put differently, the
13 observed association between GA3 and community assembly reflects selection on GA3
14 mediated by community composition, rather than GA3 mediated impact on overall
15 community assembly. Importantly, *B. fragilis* abundance is low compared to the
16 *Bacteroides* consortium as a whole and therefore, it is perhaps not surprising that we do
17 not detect other Gram-negative genera whose abundance are specifically lowered in
18 GA3+ microbiomes. However, we cannot rule out that competition from less common or
19 low abundance genera that fall below our detection limit might also select for the
20 retention of GA3 in adults. Indeed, there is evidence that *Bacteroides* spp and members
21 of the *Enterobacteriaceae* interact intimately, although T6SS-dependent interactions have
22 not yet been shown between these organisms (de Sablet et al. 2009; Buffie and Pamer
23 2013; Curtis et al. 2014; Charbonneau et al. 2016).

1 The T6SS is one of many antagonistic pathways whose operation is determined
2 by the presence or absence of polymorphic toxins and corresponding antitoxins (Aoki et
3 al. 2010; Zhang et al. 2012; Cao et al. 2016). We show here that a systematic
4 characterization and large-scale computational analysis of metagenomic data can provide
5 a means of linking the presence and abundance of these crucial factors to microbial
6 community composition. Moreover, for contact-dependent pathways such as the T6SS,
7 such analyses can provide a unique window into community biogeography. Indeed,
8 *Bacteroides* spp. are thought to occupy a crowded niche proximal to the gut mucosa and
9 our findings herein provide evidence of extensive cell–cell contacts between species of
10 the genus (Earle et al. 2015; Donaldson et al. 2016). Thus, our study offers an analytical
11 framework for more globally deciphering the forces that dictate the establishment and
12 maintenance of bacterial communities.
13

1 **Methods**

2 **Short read metagenomic datasets and genomic data**

3 Our analysis utilizes short read metagenomic data from several large-scale microbiome
4 datasets. For adult microbiomes we downloaded 147 shotgun samples from HMP
5 (Human Microbiome Project 2012), and 99 healthy human shotgun samples from
6 MetaHIT (Qin et al. 2010). Since an excessive fraction of human DNA will likely not
7 markedly impact our ability to quantify *B. fragilis* abundance, HMP samples which failed
8 QC were nonetheless included in our analysis. For infant microbiomes, we downloaded
9 300 samples from a study of development of the microbiome in the first year of life
10 (Backhed et al. 2015), 769 samples from a study of autoimmune diseases (Vatanen et al.
11 2016), 237 samples from a study of antibiotic usage (Yassour et al. 2016), and 126
12 samples from a study of the development of Type 1 Diabetes (Kostic et al. 2015). Several
13 of these datasets include multiple longitudinal samples from the same individuals, which
14 were used for temporal analysis.

15 We downloaded all available *B. fragilis* genomes from RefSeq. Sequences from 3
16 strains were found to be contaminated with contigs matching species other than *B.*
17 *fragilis*, and were discarded. A group of 8 strains appeared to be very distant in sequence
18 homology from the rest of the strains, and were also discarded. We additionally
19 downloaded from RefSeq all genomes of other *Bacteroides* species for which at least 5
20 strain genomes were available.

21

22 **Identifying species-specific marker genes**

1 We compile a list of marker genes that could be used for strain-level inference. Our
2 marker gene approach is similar to that used by MetaPhlAn (Truong et al. 2015), but rely
3 on a more stringent selection of marker genes, supporting a more robust comparison at
4 the strain level. Specifically, for our analysis, we identified a subset of the MetaPhlAn
5 marker genes that are found in the genome of *every* sequenced strain in a single copy. To
6 this end, for each of the MetaPhlAn marker genes associated with a given species, we
7 used BLASTn to find every homolog (>60% identity) in all the strains of that species. We
8 used Usearch (Edgar 2010) to cluster this larger set of genes into groups with >90%
9 identity and if a cluster with exactly one gene in each strain could be found, the marker
10 gene was included in our list.

11

12 **Identifying T6SS effector and immunity genes**

13 We sought to comprehensively catalog *Bacteroidales* T6SS E–I genes from reference
14 genomes. In *Bacteroidales*, as in other bacteria, E–I genes are encoded adjacent to the
15 genes for secreted structural proteins Hcp and VgrG. Accordingly, we manually curated
16 genes adjacent to these structural genes across all publically available *Bacteroidales*
17 genomes. Identified genes exhibited reduced GC content relative to the rest of the T6SS
18 locus or the genome as a whole, and were encoded in bi-cistrons. Putative effectors
19 always lacked characteristic signal peptides, consistent with transport via the T6SS
20 apparatus, while putative immunity genes often encoded proteins with signal peptides.
21 We used structural homology prediction (Phyre)(Kelley et al. 2015) and remote sequence
22 homology search algorithms (Hmmer)(Finn et al. 2015) to predict functions for these
23 genes, identifying many genes with functions associated with known T6SS toxin

1 effectors. As in Coyne *et al* (Coyne et al. 2016), our list included predicted cell-wall
2 degrading enzymes, lipases, and nucleases as well as putative effector domains fused to
3 either PAAR domains (DUF4280) or Hcp.

4

5 **Estimating the abundance of species-specific genes and T6SS genes in metagenomic** 6 **samples**

7 We aligned shotgun reads single end using Bowtie2 (using parameters $-a -N 1$) to the set
8 of genes of interest. Alignments with less than 97% identity, a quality score below 20, or
9 multiple hits were discarded. To quantify the abundance of each gene, the number of
10 reads aligned to this gene was normalized by the length of the gene and the total number
11 of reads in the sample.

12 The average abundance of species-specific marker genes identified above was
13 used as a proxy for the abundance of that species in the sample. We defined samples as
14 having *B. fragilis* present if at least 100 reads could be aligned to *B. fragilis*-specific
15 marker genes. When characterizing strain replacement, for which higher coverage of *B.*
16 *fragilis* genes is required, we used instead a threshold of 500 reads. Because GA1 and
17 GA2 are not restricted to a single species, and because *Bacteroidales* composition can
18 vary dramatically, we considered samples with 100 reads mapping to the GA1 and GA2
19 E-I genes.

20

21 **Nucleotide diversity calculation**

22 To estimate nucleotide diversity across species-specific marker genes, we again aligned
23 all short reads in each sample to these genes. The obtained alignments were converted

1 into a pileup using mpileup from samtools (parameters --excl-flags
2 UNMAP,QCFAIL,DUP -A -q0 -C0 -B), and finally into an allele count matrix. The first
3 and last 10 bases of each gene were discarded from the allele count matrix as we found
4 they contained many poor quality alignments. We focused on high-coverage loci only,
5 ignoring all loci where the coverage was less than 5X. If the number of high coverage
6 sites was <10% of the total length of the sequence, the sample was excluded from further
7 consideration. Variable sites were defined as those having at least 2 counts of the minor
8 allele. Nucleotide diversity was then calculated at these variant sites according to:

9
$$\pi = \frac{1}{n} \sum_i 2p_i q_i$$

10 Where n is the total length of the genes, i corresponds to the variable sites, and p and q
11 correspond to the frequency of the major and minor allele at site i .

12

13 **Inferring *B. fragilis* strains using StrainFinder**

14 To infer strain diversity in each sample, we used StrainFinder
15 (<https://github.com/cssmillie/fmt>), a previously introduced method for inferring the most
16 likely set of strains in metagenomic samples based on nucleotide variants. To this end, we
17 again aligned the short read in each sample to the set of *B. fragilis*-specific marker genes
18 identified above, and converted the alignment to a count matrix describing the number of
19 counts of each nucleotide at every position along the genes. As when calculating
20 nucleotide diversity, we discarded the first and last 10 bases and only considered sites
21 with at least 5 counts. We also discarded samples where the high coverage sites were less
22 than 10% of the total length as they resulted in poor quality trees. When running
23 StrainFinder we reduced the data to only those sites with population variability, defined,

1 as above, as sites with at least 2 counts of the minor allele. As noted above, we only
2 considered samples with a sufficient coverage on the *B. fragilis* marker genes to enable
3 robust strain inference.

4 StrainFinder determines the relative strain abundance and genotypes at variable
5 sites by considering the likelihood of the observed allele counts and using an expectation
6 maximization approach. An optimal number of strains between 1 and 10 was determined
7 using AIC. For each run of StrainFinder we used 5 independent runs of 200 expectation
8 maximization iterations and selected the best fit; these parameters yield reproducible
9 inference of strains. When analyzing temporal data with StrainFinder, we combined allele
10 counts from all samples of an individual into a single 3-dimensional matrix. Using the
11 genotypes from StrainFinder we then reconstructed strain-specific versions of each
12 marker gene and then created subsequences consisting of only the high coverage sites we
13 used in the analysis. Inferred strains were then further examined using a phylogenetic
14 analysis as described below.

15

16 **Phylogenetic analysis**

17 A phylogenetic tree of the reference *B. fragilis* strains was constructed based on their
18 species-specific marker genes. Specifically, we aligned the strains' versions of each
19 marker gene using MAFFT, concatenated the alignments of all genes, and then
20 constructed a tree using the GTRGAMMAI model from RAxML (Stamatakis 2014), as
21 has been done previously (Wexler et al. 2016).

22 To determine whether the strains inferred by StrainFinder are monophyletic, we
23 combined the sequences from the inferred strains (as determined by StrainFinder), with

1 the sequences from the available reference *B. fragilis* genomes, and recreated the strain
2 phylogeny using the same method as described above. We defined inferred strains as
3 monophyletic if their common ancestor does not have any descendent outside the set of
4 inferred strains or if the distance was less than 0.001 substitutions per site.

5

6 **Strain sequencing**

7 All human studies were conducted with the permission of the Yale Human Investigation
8 Committee. Stool samples from four healthy individuals frozen in sterile glycerol (Cullen
9 et al. 2015) were plated onto *Bacteroides* Bile Esculin or *Brucella* Blood Agar plates (BD
10 Biosciences), to select for *Bacteroidales* colonies. Single colonies were picked into
11 Mega Medium (Wu et al. 2015) and grown to stationary phase in anaerobic conditions
12 before freezing in 10% glycerol in 96-well plates. PCR was performed directly from the
13 frozen stocks using primers to amplify the V1-V4 region of the 16S rRNA gene. PCR
14 products were then sent for Sanger sequencing. Reads were converted to fastq format
15 and NCBI Blast 2.2.31+ was used to align sequences to the SILVA 123 and GreenGenes
16 2011-1 16S rRNA gene databases in order to identify *Bacteroides fragilis* colonies. To
17 verify the *B. fragilis*-positive colonies, a second round of PCR was performed using
18 primers to amplify and sequence the *gyrB* gene. Two of the four donors were confirmed
19 to have *B. fragilis*. Twenty confirmed colonies from each *B. fragilis*-positive donor were
20 then grown up to stationary phase in TYG medium under anaerobic conditions. Genomic
21 DNA was isolated using the Qiagen DNeasy Blood and Tissue Kit and prepared for
22 whole genome sequencing using the MiSeq V3 Reagent Kit. Sequencing was performed
23 in the Nickerson lab core facility in the UW Department of Genome Sciences.

1 Sequencing reads were mapped to the set of *B. fragilis*-specific marker genes to generate
2 alignments. Samples under 10x mean alignment read coverage were then discarded.
3 Consensus sequence for each remaining sample was generated using the GATK
4 FastaAlternateReferenceMaker. Subsequently, we constructed multi-alignments for all
5 the samples using MAFFT 7.237, concatenated them, and then inferred a phylogenetic
6 tree using the GTRGAMMAI model from RAxML 8.2.8 (Stamatakis 2014). All
7 sequences were deposited into NCBI SRA under BioProject ID PRJNA375094.

8

9 **Predicting E–I gene content of inferred *B. fragilis* strains**

10 We predicted the E–I gene content of an inferred *B. fragilis* strain by examining the E–I
11 content in the genome of its nearest neighbors on a phylogenetic tree that contains the
12 inferred strains from a given sample and the reference strains (as described above).
13 Specifically, for every strain identified from StrainFinder, we identified the most recent
14 ancestor that have both the inferred strain and at least one reference strain as descendants.
15 We then used the average E–I content of all reference strains descendant from this
16 ancestor as the predicted E–I content of the inferred strain. To then estimate the predicted
17 E–I content in the sample, we combined the predicted E–I content of each inferred strain
18 weighted by their relative abundance. To determine the confidence of the predicted E–I
19 content we determined the average phylogenetic distance of these reference strains to the
20 ancestor identified above.

21

22 **Classifying microbiome samples as T6SS+ vs. T6SS-**

23 For every sample, we estimated the number of reads expected to map to the *B. fragilis*

1 GA3 T6SS structural genes based on the number of reads mapped to *B. fragilis*-specific
2 marker genes in that sample and the ratio between the total length of *B. fragilis*-specific
3 marker genes and *B. fragilis* T6SS structural genes. We define samples to be T6SS+ if *B.*
4 *fragilis* was present (as defined above) and the number of reads mapped to T6SS
5 structural genes was more than 10% of the expected number (and see Figure S6A). We
6 define samples to be T6SS- if *B. fragilis* was present and the number of reads mapped to
7 T6SS structural genes was less than 10% of the expected number.

8

9 **Community composition analysis**

10 To obtain independent estimate community composition in each sample, we
11 downloaded the v35 16S OTU abundance table for human gut microbiomes from HMP
12 (ftp://public-ftp.hmpdacc.org/HMQCP/otu_table_psn_v35.txt.gz), summed the counts
13 from all OTUs in the same genus, and calculated the relative abundance of each genus.
14 Importantly, because 16S sequencing depth is independent from the depth of shotgun
15 samples used to determine T6SS+ vs. T6SS- classification, using these 16S-based data
16 allows us to compare T6SS presence with community taxonomic profiles without
17 potential coverage-related biases. Samples were classified into T6SS+ and T6SS- as
18 described above. GA1 and GA2 lack uniquely identifying structural genes so we defined
19 T6SS+ vs. T6SS- as samples with vs. without 100 counts mapping to the GA1 or GA2 E-
20 I genes respectively. The distance between samples was defined by the Bray-Curtis
21 distance at the genus level, and significance of separation between T6SS+ and T6SS-
22 samples was evaluated using PERMANOVA. For the subset of genera whose average
23 abundance across samples was > 0.1%, we used a Wilcoxon rank sum test to compare

1 their abundance in T6SS+ vs. T6SS- samples using a 5% FDR.

2

3 **Acknowledgements**

4 We thank Eric Alm and Chris Smillie for sharing StrainFinder code and for their support
5 in running it. We also thank UW Genome Sciences ITS for high-performance computing
6 resources. We are grateful to S. Brook Peterson for careful review of the manuscript, and
7 to members of the Borenstein and Mougous laboratories for helpful discussions. This
8 work was supported by National Institutes of Health grants GM118159 (to ALG),
9 AI080609 (to JDM), and New Innovator Award DP2AT00780201 (to EB), the Pew
10 Scholars Program (ALG), and the Burroughs Wellcome Fund (ALG and JDM). AJV was
11 supported by a postdoctoral fellowship from the Natural Sciences and Engineering
12 Research Council of Canada. BDR was supported by a Simons Foundation-sponsored
13 Life Sciences Research Foundation postdoctoral fellowship.

14

15

1 **References**

- 2
- 3 Aoki SK, Diner EJ, de Roodenbeke CT, Burgess BR, Poole SJ, Braaten BA, Jones AM,
4 Webb JS, Hayes CS, Cotter PA et al. 2010. A widespread family of polymorphic
5 contact-dependent toxin delivery systems in bacteria. *Nature* **468**: 439-442.
- 6 Backhed F, Roswall J, Peng Y, Feng Q, Jia H, Kovatcheva-Datchary P, Li Y, Xia Y, Xie
7 H, Zhong H et al. 2015. Dynamics and Stabilization of the Human Gut
8 Microbiome during the First Year of Life. *Cell Host Microbe* **17**: 690-703.
- 9 Buffie CG, Pamer EG. 2013. Microbiota-mediated colonization resistance against
10 intestinal pathogens. *Nat Rev Immunol* **13**: 790-801.
- 11 Bulgarelli D, Garrido-Oter R, Munch PC, Weiman A, Droge J, Pan Y, McHardy AC,
12 Schulze-Lefert P. 2015. Structure and function of the bacterial root microbiota in
13 wild and domesticated barley. *Cell Host Microbe* **17**: 392-403.
- 14 Cao Z, Casabona MG, Kneuper H, Chalmers JD, Palmer T. 2016. The type VII secretion
15 system of *Staphylococcus aureus* secretes a nuclease toxin that targets competitor
16 bacteria. *Nat Microbiol* **2**: 16183.
- 17 Charbonneau MR, O'Donnell D, Blanton LV, Totten SM, Davis JC, Barratt MJ, Cheng J,
18 Guruge J, Talcott M, Bain JR et al. 2016. Sialylated Milk Oligosaccharides
19 Promote Microbiota-Dependent Growth in Models of Infant Undernutrition. *Cell*
20 **164**: 859-871.
- 21 Chatzidaki-Livanis M, Coyne MJ, Comstock LE. 2014. An antimicrobial protein of the
22 gut symbiont *Bacteroides fragilis* with a MACPF domain of host immune
23 proteins. *Mol Microbiol* **94**: 1361-1374.
- 24 Chatzidaki-Livanis M, Geva-Zatorsky N, Comstock LE. 2016. *Bacteroides fragilis* type
25 VI secretion systems use novel effector and immunity proteins to antagonize
26 human gut *Bacteroidales* species. *Proc Natl Acad Sci U S A* **113**: 3627-3632.
- 27 Coyne MJ, Roelofs KG, Comstock LE. 2016. Type VI secretion systems of human gut
28 *Bacteroidales* segregate into three genetic architectures, two of which are
29 contained on mobile genetic elements. *BMC Genomics* **17**: 58.
- 30 Coyne MJ, Zitomersky NL, McGuire AM, Earl AM, Comstock LE. 2014. Evidence of
31 extensive DNA transfer between *bacteroidales* species within the human gut.
32 *MBio* **5**: e01305-01314.
- 33 Cullen TW, Schofield WB, Barry NA, Putnam EE, Rundell EA, Trent MS, Degnan PH,
34 Booth CJ, Yu H, Goodman AL. 2015. Gut microbiota. Antimicrobial peptide
35 resistance mediates resilience of prominent gut commensals during inflammation.
36 *Science* **347**: 170-175.

- 1 Curtis MM, Hu Z, Klimko C, Narayanan S, Deberardinis R, Sperandio V. 2014. The gut
2 commensal *Bacteroides thetaiotaomicron* exacerbates enteric infection through
3 modification of the metabolic landscape. *Cell Host Microbe* **16**: 759-769.
- 4 de Sablet T, Chassard C, Bernalier-Donadille A, Vareille M, Gobert AP, Martin C. 2009.
5 Human microbiota-secreted factors inhibit shiga toxin synthesis by
6 enterohemorrhagic *Escherichia coli* O157:H7. *Infect Immun* **77**: 783-790.
- 7 Donaldson GP, Lee SM, Mazmanian SK. 2016. Gut biogeography of the bacterial
8 microbiota. *Nat Rev Microbiol* **14**: 20-32.
- 9 Dong TG, Ho BT, Yoder-Himes DR, Mekalanos JJ. 2013. Identification of T6SS-
10 dependent effector and immunity proteins by Tn-seq in *Vibrio cholerae*. *Proc Natl*
11 *Acad Sci U S A* **110**: 2623-2628.
- 12 Earle KA, Billings G, Sigal M, Lichtman JS, Hansson GC, Elias JE, Amieva MR, Huang
13 KC, Sonnenburg JL. 2015. Quantitative Imaging of Gut Microbiota Spatial
14 Organization. *Cell Host Microbe* **18**: 478-488.
- 15 Edgar RC. 2010. Search and clustering orders of magnitude faster than BLAST.
16 *Bioinformatics* **26**: 2460-2461.
- 17 Egan F, Reen FJ, O'Gara F. 2015. The distribution and diversity in metagenomic datasets
18 reveal niche specialization. *Environ Microbiol Rep* **7**: 194-203.
- 19 Faith JJ, Guruge JL, Charbonneau M, Subramanian S, Seedorf H, Goodman AL,
20 Clemente JC, Knight R, Heath AC, Leibel RL et al. 2013. The long-term stability
21 of the human gut microbiota. *Science* **341**: 1237439.
- 22 Finn RD, Clements J, Arndt W, Miller BL, Wheeler TJ, Schreiber F, Bateman A, Eddy
23 SR. 2015. HMMER web server: 2015 update. *Nucleic Acids Res* **43**: W30-38.
- 24 Hecht AL, Casterline BW, Earley ZM, Goo YA, Goodlett DR, Bubeck Wardenburg J.
25 2016. Strain competition restricts colonization of an enteric pathogen and prevents
26 colitis. *EMBO Rep* **17**: 1281-1291.
- 27 Hibbing ME, Fuqua C, Parsek MR, Peterson SB. 2010. Bacterial competition: surviving
28 and thriving in the microbial jungle. *Nat Rev Microbiol* **8**: 15-25.
- 29 Hood RD, Singh P, Hsu F, Guvener T, Carl MA, Trinidad RR, Silverman JM, Ohlson
30 BB, Hicks KG, Plemel RL et al. 2010. A type VI secretion system of
31 *Pseudomonas aeruginosa* targets a toxin to bacteria. *Cell Host Microbe* **7**: 25-37.
- 32 Human Microbiome Project C. 2012. Structure, function and diversity of the healthy
33 human microbiome. *Nature* **486**: 207-214.
- 34 Kau AL, Ahern PP, Griffin NW, Goodman AL, Gordon JI. 2011. Human nutrition, the
35 gut microbiome and the immune system. *Nature* **474**: 327-336.

- 1 Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015. The Phyre2 web portal
2 for protein modeling, prediction and analysis. *Nat Protoc* **10**: 845-858.
- 3 Kostic AD, Gevers D, Siljander H, Vatanen T, Hyotylainen T, Hamalainen AM, Peet A,
4 Tillmann V, Poho P, Mattila I et al. 2015. The dynamics of the human infant gut
5 microbiome in development and in progression toward type 1 diabetes. *Cell Host*
6 *Microbe* **17**: 260-273.
- 7 Ley RE, Peterson DA, Gordon JI. 2006. Ecological and evolutionary forces shaping
8 microbial diversity in the human intestine. *Cell* **124**: 837-848.
- 9 Morris JJ, Lenski RE, Zinser ER. 2012. The Black Queen Hypothesis: evolution of
10 dependencies through adaptive gene loss. *MBio* **3**.
- 11 Nayfach S, Rodriguez-Mueller B, Garud N, Pollard KS. 2016. An integrated
12 metagenomics pipeline for strain profiling reveals novel patterns of bacterial
13 transmission and biogeography. *Genome Research* doi:10.1101/gr.201863.115.
- 14 Prosser JI, Bohannan BJ, Curtis TP, Ellis RJ, Firestone MK, Freckleton RP, Green JL,
15 Green LE, Killham K, Lennon JJ et al. 2007. The role of ecological theory in
16 microbial ecology. *Nat Rev Microbiol* **5**: 384-392.
- 17 Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons N,
18 Levenez F, Yamada T et al. 2010. A human gut microbial gene catalogue
19 established by metagenomic sequencing. *Nature* **464**: 59-65.
- 20 Riley MA, Wertz JE. 2002. Bacteriocins: evolution, ecology, and application. *Annu Rev*
21 *Microbiol* **56**: 117-137.
- 22 Roelofs KG, Coyne MJ, Gentyala RR, Chatzidaki-Livanis M, Comstock LE. 2016.
23 Bacteroidales Secreted Antimicrobial Proteins Target Surface Molecules
24 Necessary for Gut Colonization and Mediate Competition In Vivo. *MBio* **7**.
- 25 Russell AB, Peterson SB, Mougous JD. 2014a. Type VI secretion system effectors:
26 poisons with a purpose. *Nat Rev Microbiol* **12**: 137-148.
- 27 Russell AB, Wexler AG, Harding BN, Whitney JC, Bohn AJ, Goo YA, Tran BQ, Barry
28 NA, Zheng H, Peterson SB et al. 2014b. A type VI secretion-related pathway in
29 Bacteroidetes mediates interbacterial antagonism. *Cell Host Microbe* **16**: 227-236.
- 30 Sana TG, Flaughnatti N, Lugo KA, Lam LH, Jacobson A, Baylot V, Durand E, Journet L,
31 Cascales E, Monack DM. 2016. Salmonella Typhimurium utilizes a T6SS-
32 mediated antibacterial weapon to establish in the host gut. *Proc Natl Acad Sci U S*
33 *A* **113**: E5044-5051.
- 34 Schwarz S, Hood RD, Mougous JD. 2010. What is type VI secretion doing in all those
35 bugs? *Trends Microbiol* **18**: 531-537.

- 1 Sharon I, Morowitz MJ, Thomas BC, Costello EK, Relman DA, Banfield JF. 2013. Time
2 series community genomics analysis reveals rapid shifts in bacterial species,
3 strains, and phage during infant gut colonization. *Genome Res* **23**: 111-120.
- 4 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis
5 of large phylogenies. *Bioinformatics* **30**: 1312-1313.
- 6 Trosvik P, de Muinck EJ. 2015. Ecology of bacteria in the human gastrointestinal tract--
7 identification of keystone and foundation taxa. *Microbiome* **3**: 44.
- 8 Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A,
9 Huttenhower C, Segata N. 2015. MetaPhlan2 for enhanced metagenomic
10 taxonomic profiling. *Nat Methods* **12**: 902-903.
- 11 Unterwiesing D, Miyata ST, Bachmann V, Brooks TM, Mullins T, Kostiuk B, Provenzano
12 D, Pukatzki S. 2014. The *Vibrio cholerae* type VI secretion system employs
13 diverse effector modules for intraspecific competition. *Nat Commun* **5**: 3549.
- 14 Vatanen T, Kostic AD, d'Hennezel E, Siljander H, Franzosa EA, Yassour M, Kolde R,
15 Vlamakis H, Arthur TD, Hamalainen AM et al. 2016. Variation in Microbiome
16 LPS Immunogenicity Contributes to Autoimmunity in Humans. *Cell* **165**: 842-
17 853.
- 18 Walter J, Ley R. 2011. The human gut microbiome: ecology and recent evolutionary
19 changes. *Annu Rev Microbiol* **65**: 411-429.
- 20 Wexler AG, Bao Y, Whitney JC, Bobay LM, Xavier JB, Schofield WB, Barry NA,
21 Russell AB, Tran BQ, Goo YA et al. 2016. Human symbionts inject and
22 neutralize antibacterial toxins to persist in the gut. *Proc Natl Acad Sci U S A* **113**:
23 3639-3644.
- 24 Whitney JC, Quentin D, Sawai S, LeRoux M, Harding BN, Ledvina HE, Tran BQ,
25 Robinson H, Goo YA, Goodlett DR et al. 2015. An interbacterial NAD(P)(+)
26 glycohydrolase toxin requires elongation factor Tu for delivery to target cells.
27 *Cell* **163**: 607-619.
- 28 Wu M, McNulty NP, Rodionov DA, Khoroshkin MS, Griffin NW, Cheng J, Latreille P,
29 Kerstetter RA, Terrapon N, Henrissat B et al. 2015. Genetic determinants of in
30 vivo fitness and diet responsiveness in multiple human gut *Bacteroides*. *Science*
31 **350**: aac5992.
- 32 Yassour M, Vatanen T, Siljander H, Hamalainen AM, Harkonen T, Ryhanen SJ,
33 Franzosa EA, Vlamakis H, Huttenhower C, Gevers D et al. 2016. Natural history
34 of the infant gut microbiome and impact of antibiotic treatment on bacterial strain
35 diversity and stability. *Sci Transl Med* **8**: 343ra381.
- 36 Zhang D, de Souza RF, Anantharaman V, Iyer LM, Aravind L. 2012. Polymorphic toxin
37 systems: Comprehensive characterization of trafficking modes, processing,

1 mechanisms of action, immunity and ecology using comparative genomics. *Biol*
2 *Direct* **7**: 18.
3
4

1 **Supporting Figure Legends**

2

3 **Figure S1 | Identification of unique *B. fragilis* GA3 effectors.**

4 **A.** Schematic of the GA3 locus from *B. fragilis* NCTC 9343. GC and AT content plotted
5 below locus. **B.** GC content of indicated regions from 10 representative *B. fragilis* strains.
6 **C.** Unique *B. fragilis* GA3 E–I pairs from cassette 1 and cassette 2 used in subsequent
7 metagenomics analyses.

8

9 **Figure S2 | Minimal relative error in effector abundance assuming that the T6SS is**
10 **encoded by a single species.**

11 The relative error is defined as the difference between the average abundance of detectable
12 effector genes in a sample and the expected abundance of these genes assuming they are
13 encoded by a given species. For each sample, the minimal relative error (across all possible
14 species) is plotted and samples are ordered by the magnitude of the minimal relative error.
15 Only samples in which at least 100 reads mapped to the E–I genes of a given subsystem
16 are included. The color of each point represents the species for which the minimal relative
17 error was obtained.

18

19 **Figure S3 | The co-occurrence of GA3 effector genes in metagenomes and genomes.**

20 Each cell in the heatmap, a_{ij} , denotes the probability that effector gene i is detected in a
21 sample **A.** or encoded in a genome **B.** given that effector gene j was detected/encoded. The
22 barplots on the top and right of each heatmap illustrate the number of metagenomes or
23 genomes in which each effector gene was detected.

24

1 **Figure S4 | *Bacteroides* nucleotide and strain diversity in the gut microbiome of adults.**

2 **A.** The nucleotide diversity for different *Bacteroides* spp. in adult samples from the HMP
3 and MetaHIT studies. Nucleotide diversity was calculated based on population variants in
4 species-specific marker genes (Methods). Only species with at least 5 genomes in RefSeq
5 were considered. **B.** A phylogenetic tree linking previously sequenced *B. fragilis* reference
6 genomes with sequenced colonies from two individuals (in red and blue). The effector
7 genes encoded by each reference genome and the new sequenced genomes from stool are
8 represented by colored squares as in Figure 1A.

9

10 **Figure S5 | Measured and expected E–I gene abundances in HMP and MetaHIT**
11 **samples.**

12 Measured abundances (hollow circles) are based on short read mapping to effector genes.
13 Expected abundances (filled circles) are based on the reference strains phylogenetically
14 closest to the inferred strain. Point size is scaled to the calculated copy number of each
15 effector. The barplot on the left shows the phylogenetic distance between the inferred strain
16 and its nearest reference strain in the phylogenetic tree.

17

18 **Figure S6 | The prevalence of T6SS- samples in infants and adults.**

19 **A.** The relationship between the number of reads expected and measured to map to the
20 GA3 T6SS structural genes, from each adult (red) and infant (blue) gut microbiome sample.
21 The expected number is based on the number of reads mapping to *B. fragilis*-specific
22 marker genes, normalized by gene lengths. The dotted line represents the cutoff used for
23 determining that *B. fragilis* is present in a sample. The solid line represents an observed

1 number of reads that is 10% of expected, and was used to distinguish T6SS+ from T6SS-
2 samples. As evident by this plot, T6SS+ and T6SS- samples can be clearly defined.
3 Different shapes correspond to different datasets, with adult samples colored in red and
4 infant samples in blue. **B.** The percentage of individuals of those harboring *B. fragilis*, that
5 lack the GA3 T6SS across different adult (red) and infant (blue) datasets. Individuals that
6 were not consistent in terms of T6SS+/- classification across different time points were not
7 included.

8

9 **Figure S7 | The abundance of *B. fragilis* GA3 E–I genes in infant microbiomes.**

10 Details and format are as in Figure 1D.

11

12 **Figure S8 | All E–I turnover and strain replacement events detected in adult and**
13 **infant microbiomes.**

14 Details and format are as in Figure 3C-E (top plots). Each pair of plots is labeled with the
15 individual code and the dataset it comes from.

16

17

18

Figure S1

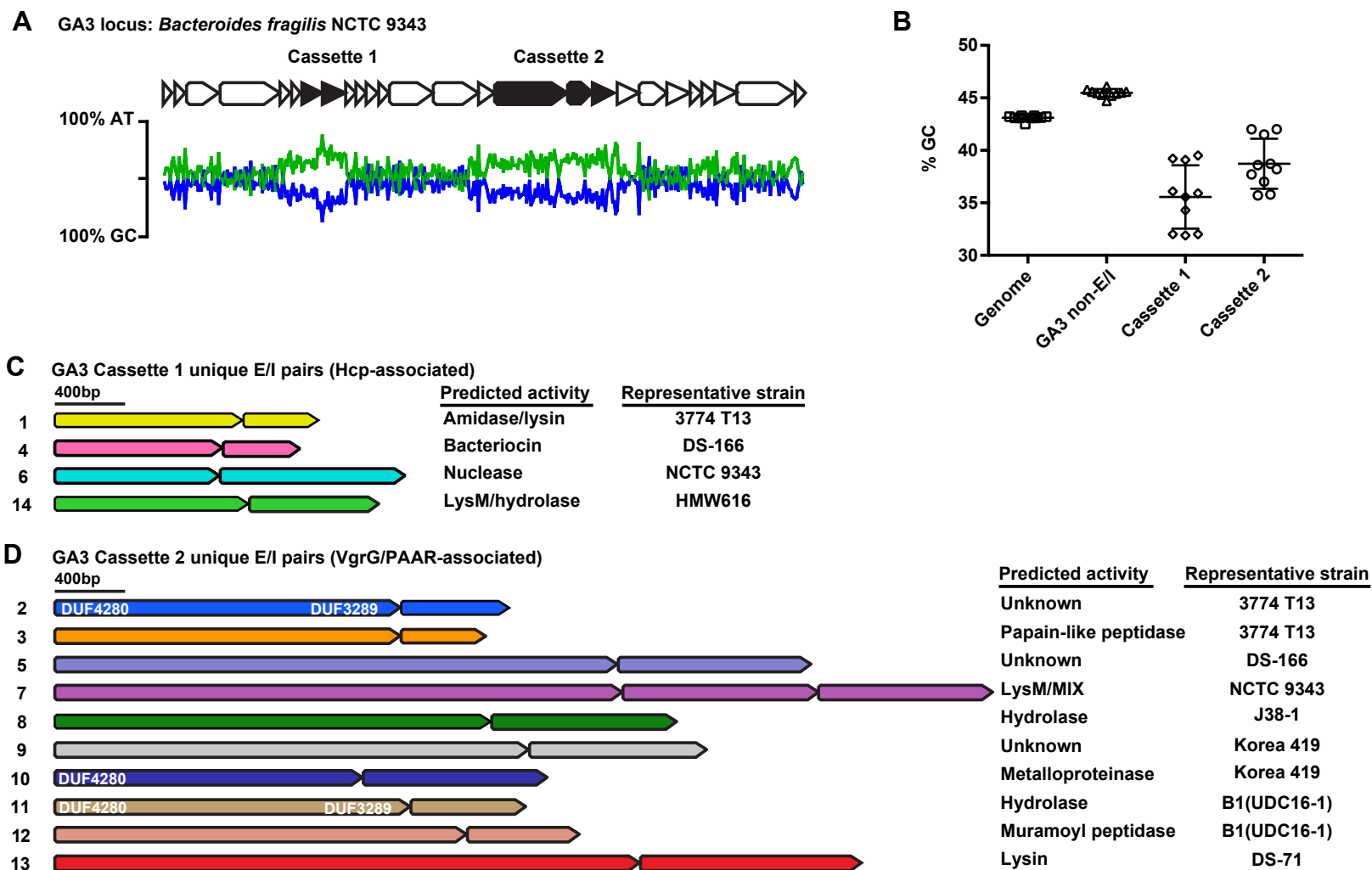


Figure S2

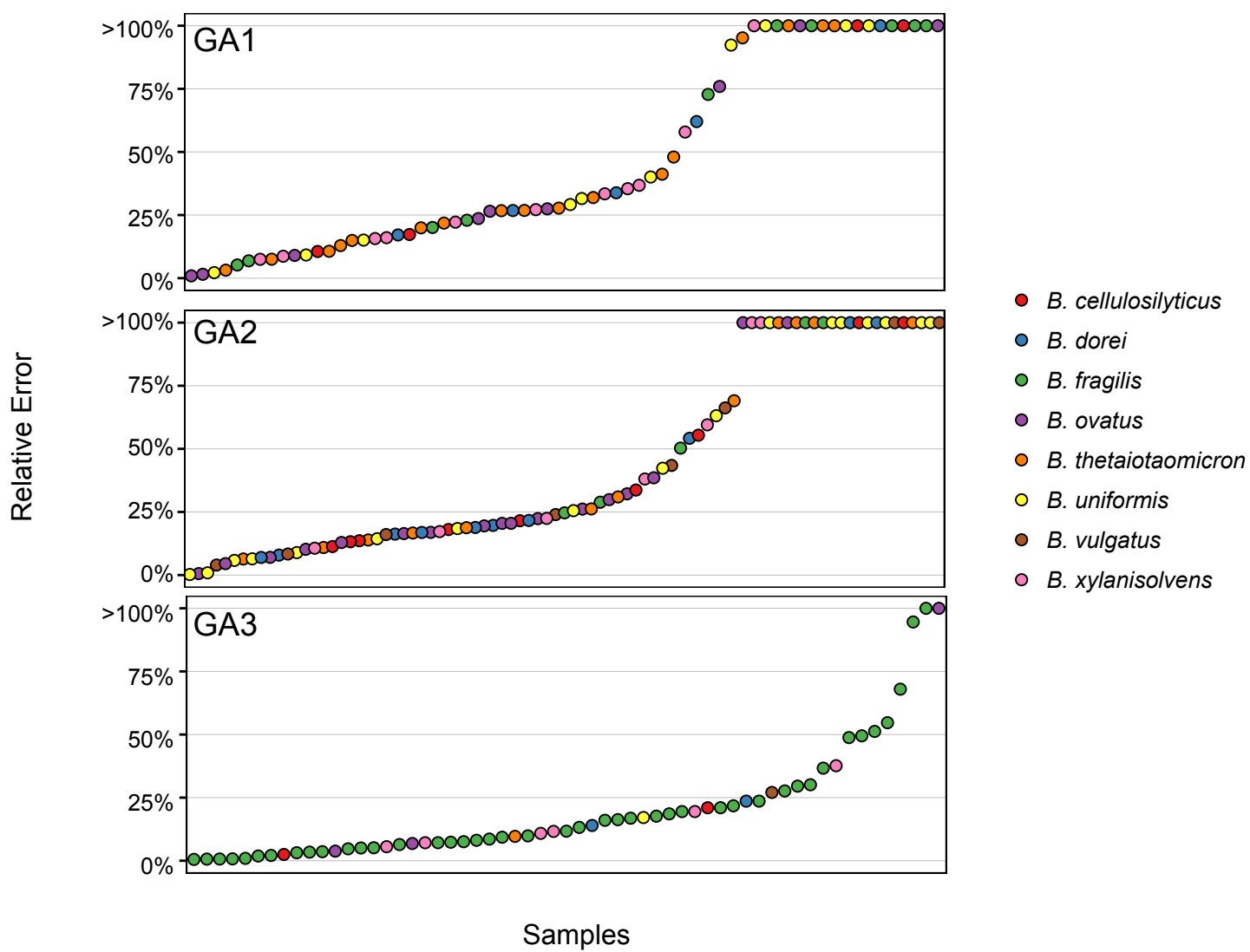
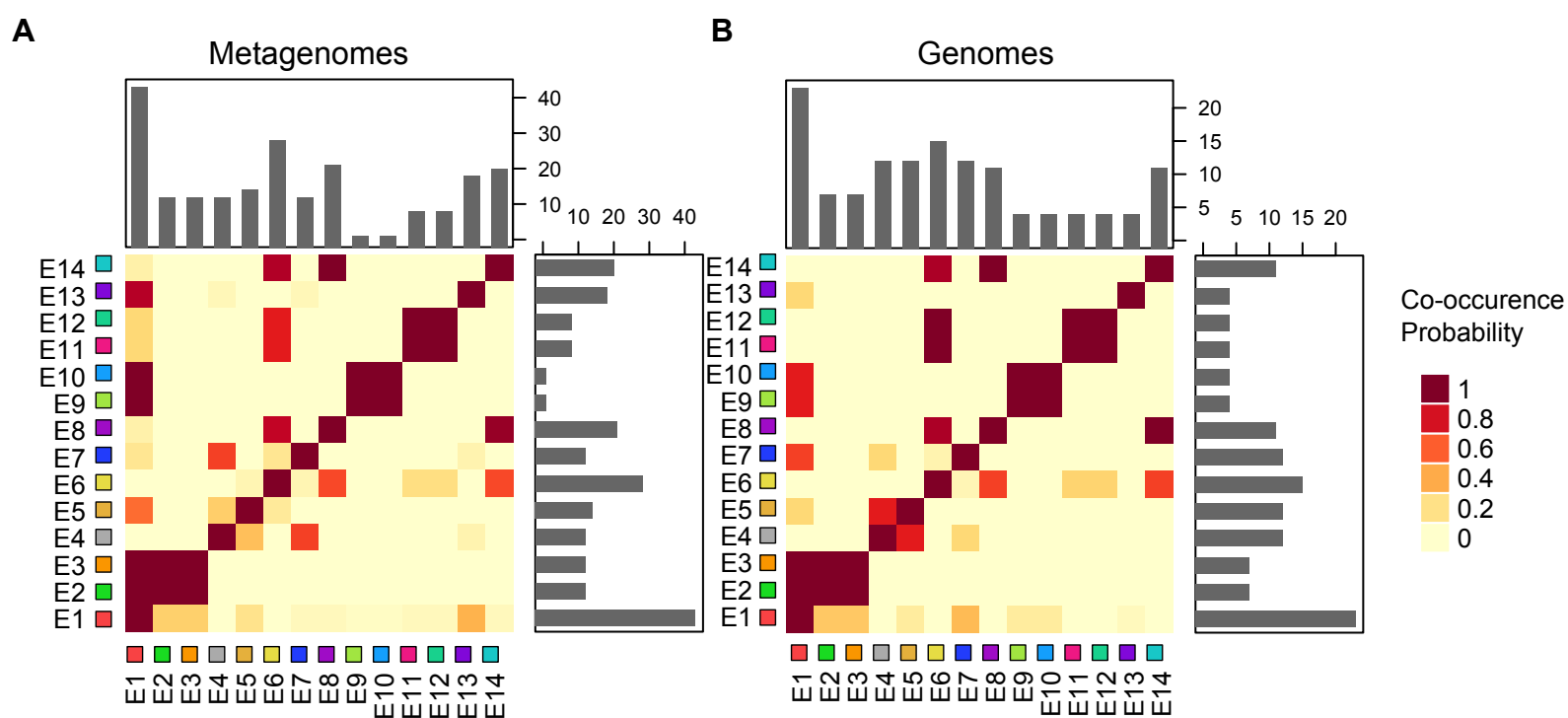
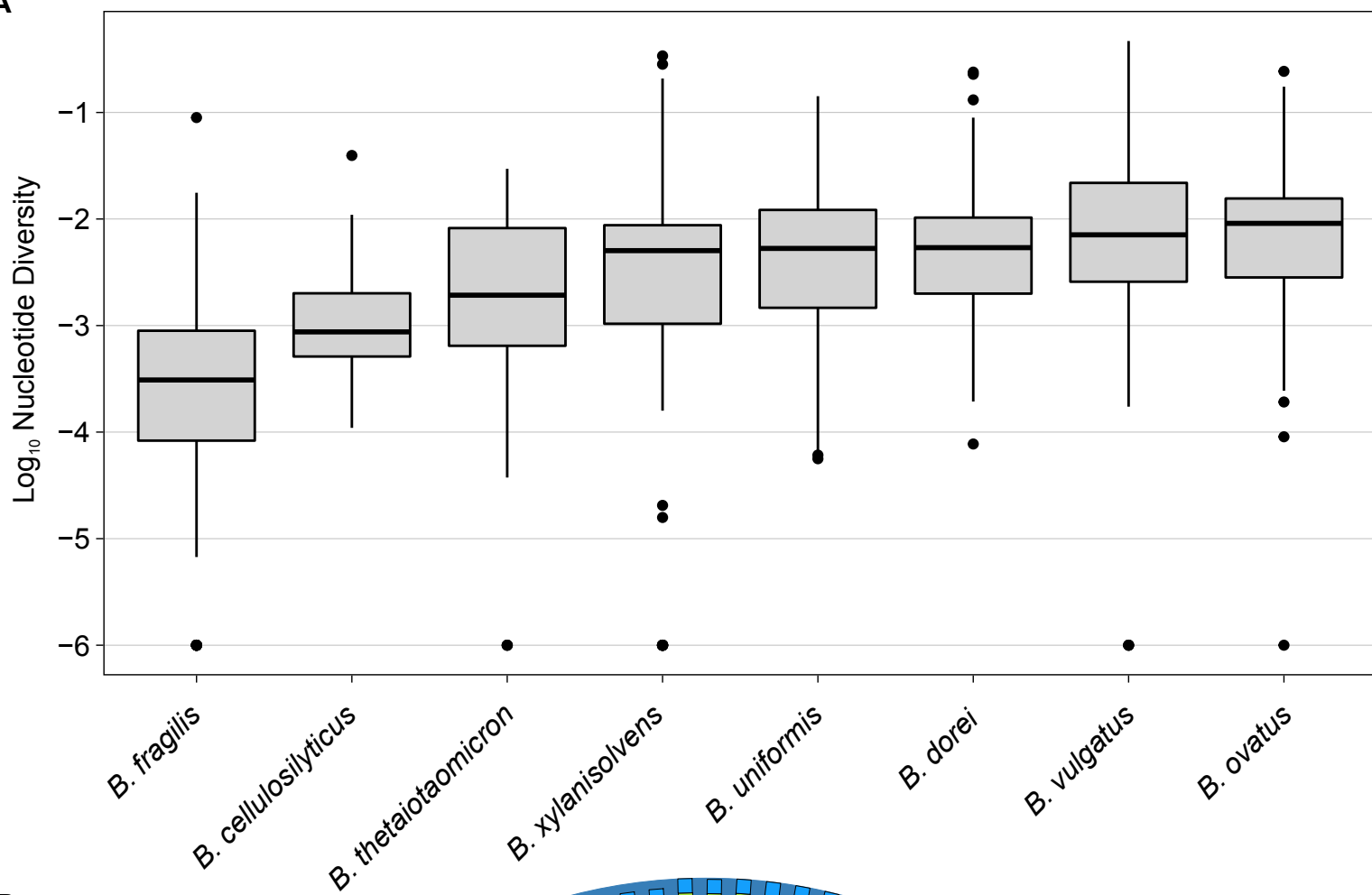


Figure S3



A



B

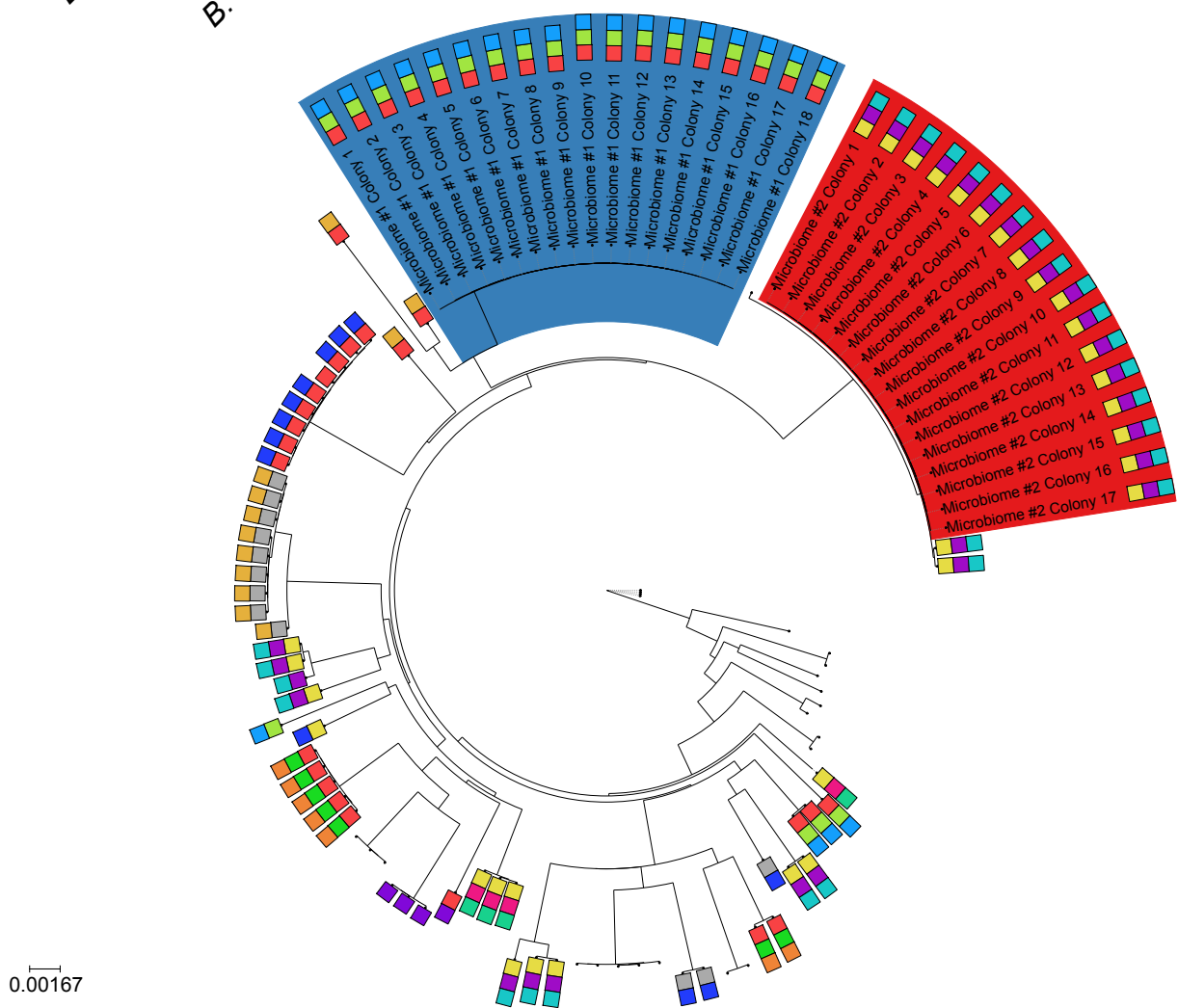


Figure S5

bioRxiv preprint doi: <https://doi.org/10.1101/134874>; this version posted May 8, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

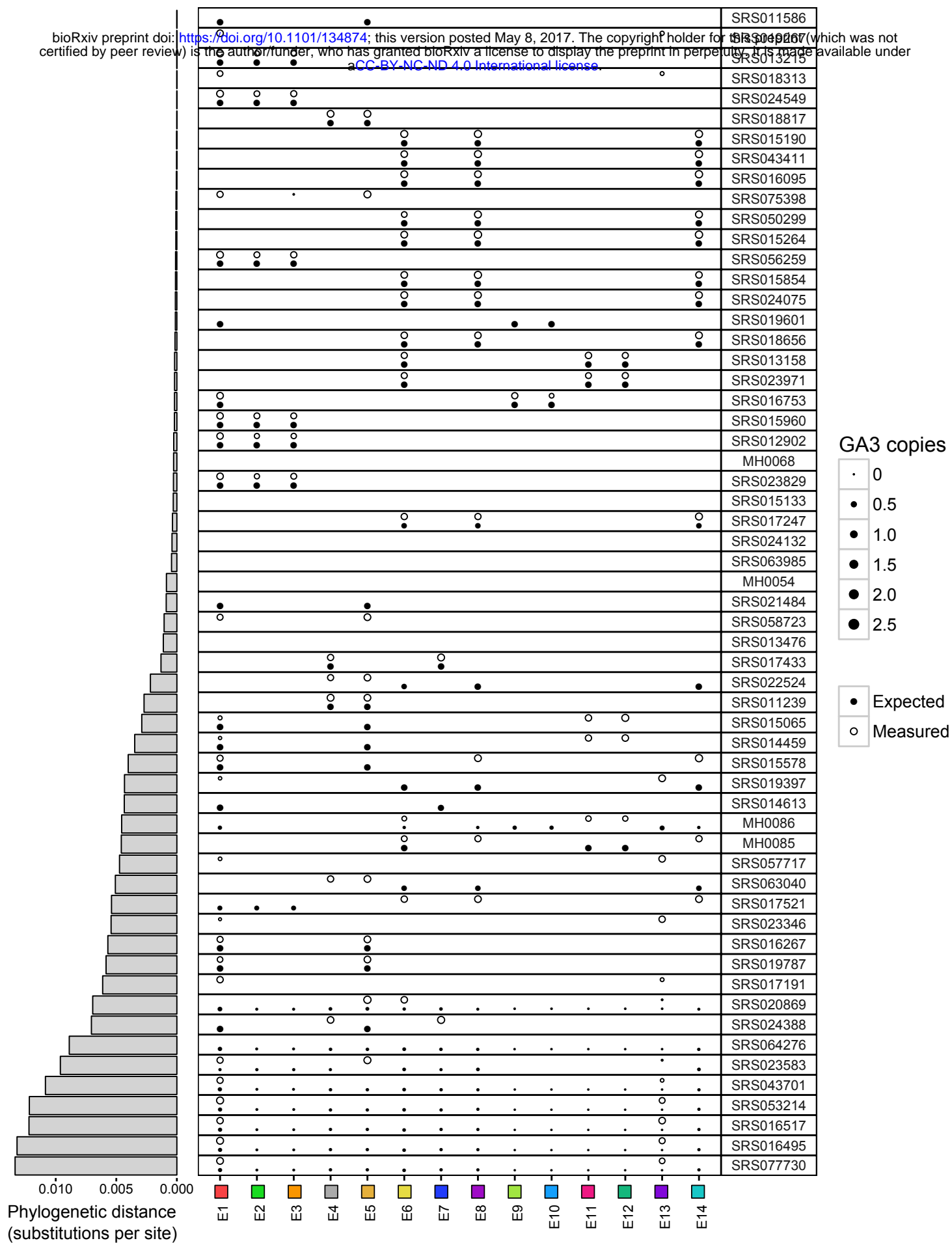
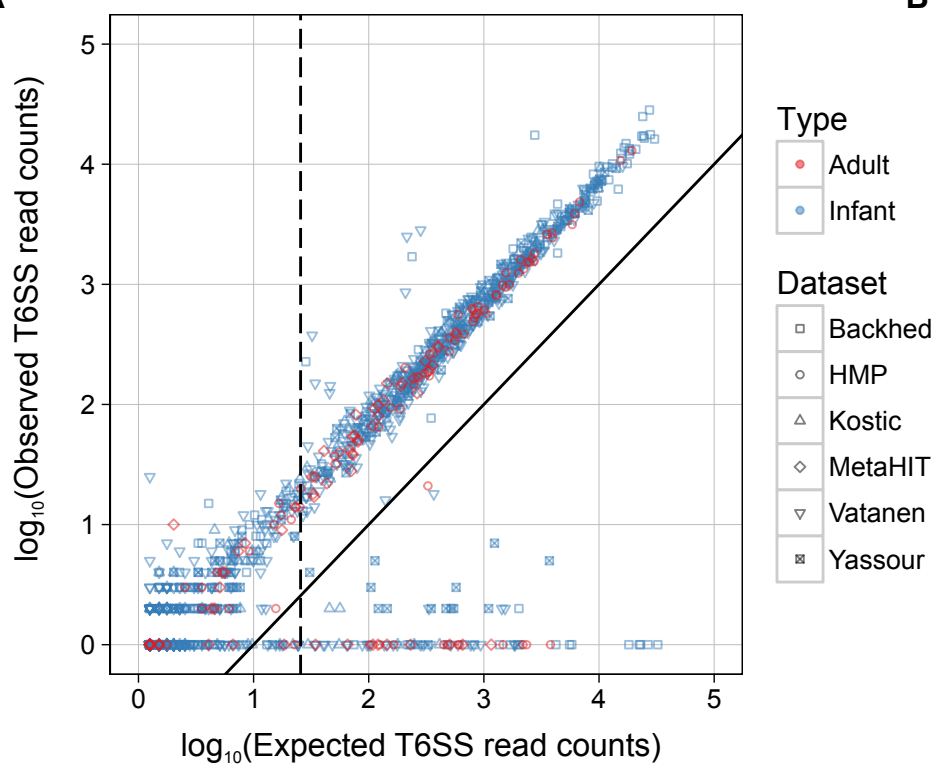


Figure S6

A



B

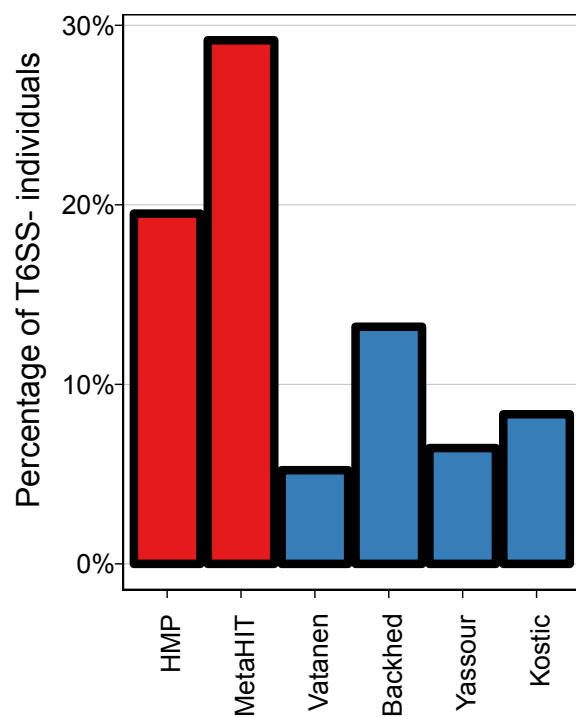


Figure S7

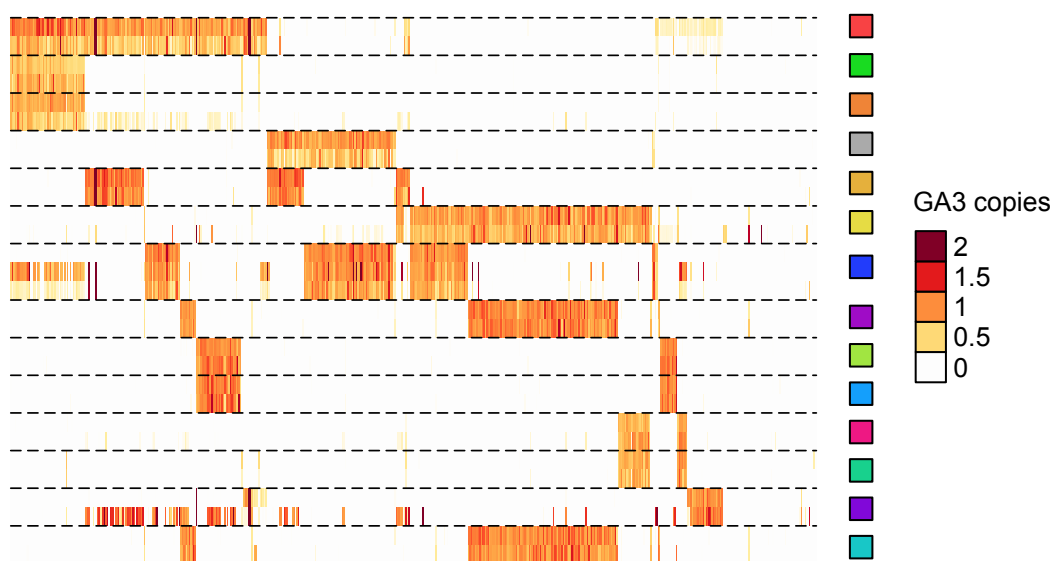


Figure S8

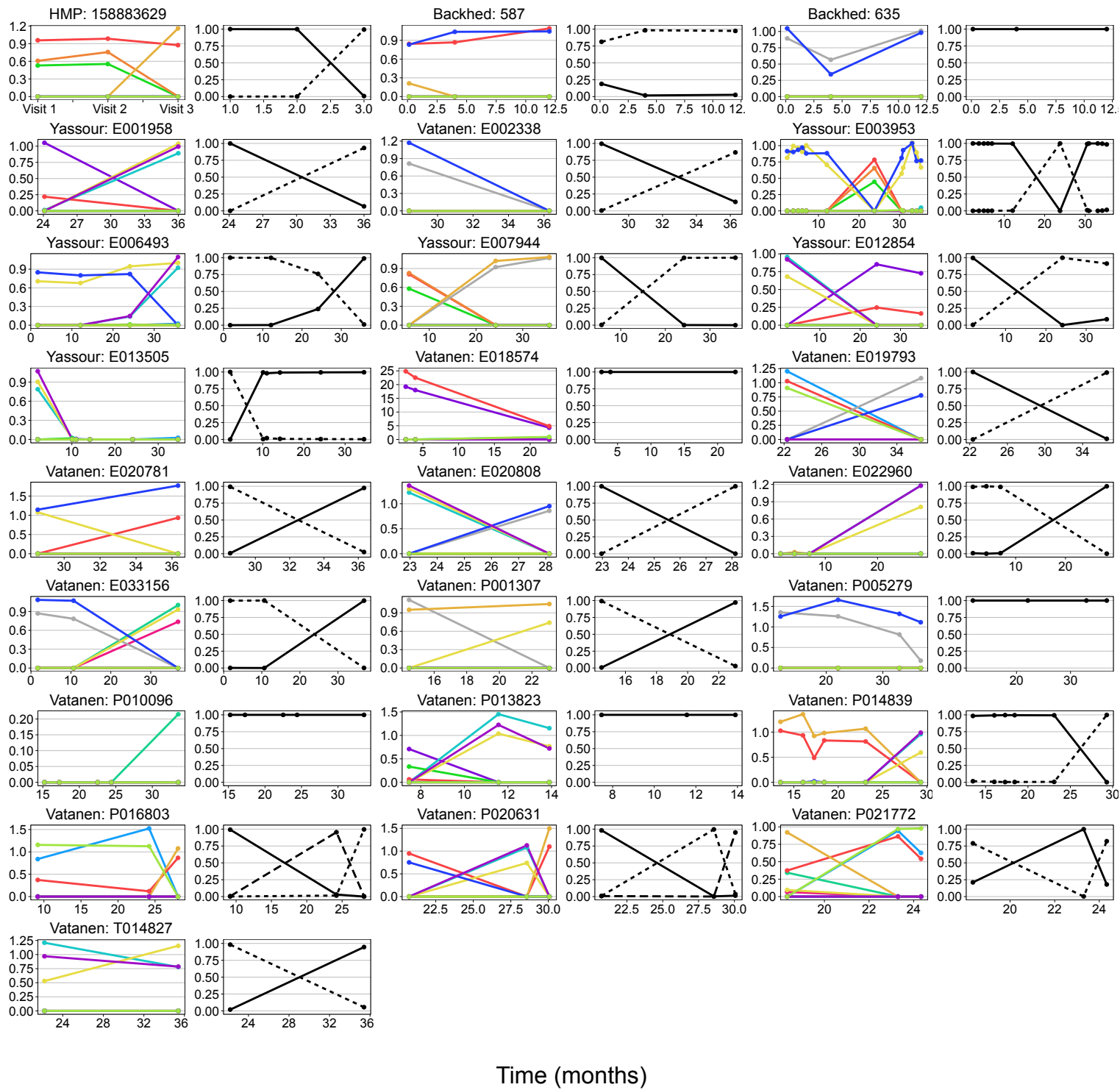


Table S1. Genus-level differential abundance in T6SS+ vs. T6SS- samples.

Genus	Mean Abundance T6SS+ (%)	Mean Abundance T6SS- (%)	P value	FDR Adjusted P value	Passes FDR
<i>Faecalibacterium</i>	3.8	10.6	8.3E-04	1.7E-02	TRUE
<i>Bacteroides</i>	57.9	37.5	3.1E-03	1.9E-02	TRUE
<i>Ruminococcus</i>	1.5	3.5	3.1E-03	1.9E-02	TRUE
<i>Oscillospira</i>	1.8	3.5	3.6E-03	1.9E-02	TRUE
<i>Eubacterium</i>	0.3	0.4	1.7E-02	7.2E-02	FALSE
<i>Odoribacter</i>	0.4	0.8	2.8E-02	9.8E-02	FALSE
<i>Subdoligranulum</i>	1.9	3.3	3.5E-02	1.1E-01	FALSE
<i>Sutterella</i>	2.4	0.1	7.6E-02	2.0E-01	FALSE
<i>Dialister</i>	1.2	0.7	1.1E-01	2.5E-01	FALSE
<i>Clostridium</i>	1.5	0.9	1.6E-01	3.2E-01	FALSE
<i>Alistipes</i>	5.7	6.7	1.8E-01	3.2E-01	FALSE
<i>Coprococcus</i>	0.5	0.6	1.8E-01	3.2E-01	FALSE
<i>Akkermansia</i>	0.1	0.3	2.2E-01	3.5E-01	FALSE
<i>Lachnospira</i>	0.7	0.7	2.6E-01	3.9E-01	FALSE
<i>Parabacteroides</i>	4.5	2.1	3.0E-01	4.2E-01	FALSE
<i>Roseburia</i>	1.5	2.7	4.5E-01	5.7E-01	FALSE
<i>Blautia</i>	0.9	0.5	4.6E-01	5.7E-01	FALSE
<i>Megamonas</i>	0.2	0.0	5.3E-01	6.2E-01	FALSE
<i>Prevotella</i>	0.1	0.3	6.0E-01	6.6E-01	FALSE
<i>Phascolarctobacterium</i>	0.8	0.8	8.0E-01	8.4E-01	FALSE
<i>Escherichia</i>	0.1	0.1	8.5E-01	8.5E-01	FALSE