

# Microbial-enrichment method enables high-throughput metagenomic characterization from host-rich samples

Received: 29 November 2022

Accepted: 27 August 2023

Published online: 12 October 2023

 Check for updates

Natalie J. Wu-Woods<sup>1,10</sup>, Jacob T. Barlow<sup>1,10</sup>, Florian Trigodet<sup>2</sup>,  
Dustin G. Shaw<sup>2,3,4</sup>, Anna E. Romano<sup>5</sup>, Bana Jabri<sup>2,3,4</sup>, A. Murat Eren<sup>6,7,8,9</sup>  
& Rustem F. Ismagilov<sup>1,5</sup> ✉

Host–microbe interactions have been linked to health and disease states through the use of microbial taxonomic profiling, mostly via 16S ribosomal RNA gene sequencing. However, many mechanistic insights remain elusive, in part because studying the genomes of microbes associated with mammalian tissue is difficult due to the high ratio of host to microbial DNA in such samples. Here we describe a microbial-enrichment method (MEM), which we demonstrate on a wide range of sample types, including saliva, stool, intestinal scrapings, and intestinal mucosal biopsies. MEM enabled high-throughput characterization of microbial metagenomes from human intestinal biopsies by reducing host DNA more than 1,000-fold with minimal microbial community changes (roughly 90% of taxa had no significant differences between MEM-treated and untreated control groups). Shotgun sequencing of MEM-treated human intestinal biopsies enabled characterization of both high- and low-abundance microbial taxa, pathways and genes longitudinally along the gastrointestinal tract. We report the construction of metagenome-assembled genomes directly from human intestinal biopsies for bacteria and archaea at relative abundances as low as 1%. Analysis of metagenome-assembled genomes reveals distinct subpopulation structures between the small and large intestine for some taxa. MEM opens a path for the microbiome field to acquire deeper insights into host–microbe interactions by enabling in-depth characterization of host-tissue-associated microbial communities.

The mucosal microbiota of the intestine has been implicated in a wide range of health conditions<sup>1</sup> including cancer<sup>2,3</sup>, inflammatory bowel disease (IBD)<sup>4–7</sup> and celiac disease<sup>8,9</sup>. Most microbiome studies use fecal samples to infer the gastrointestinal (GI) microbiota due to its

ease of access, despite microbes in feces and intestinal biopsies having distinct ecological niches<sup>10–13</sup>.

Most microbiome studies sequence 16S ribosomal RNA (rRNA) gene amplicons<sup>14</sup>, enabling detailed descriptions of taxonomic

<sup>1</sup>Biology and Bioengineering, California Institute of Technology, Pasadena, CA, USA. <sup>2</sup>Department of Medicine, The University of Chicago, Chicago, IL, USA. <sup>3</sup>Committee on Immunology, The University of Chicago, Chicago, IL, USA. <sup>4</sup>Department of Pathology, The University of Chicago, Chicago, IL, USA. <sup>5</sup>Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, USA. <sup>6</sup>Bay Paul Center, Marine Biological Laboratory, Woods Hole, MA, USA. <sup>7</sup>Institute for Chemistry and Biology of the Marine Environment, University of Oldenburg, Oldenburg, Germany. <sup>8</sup>Alfred-Wegener-Institute for Marine and Polar Research, Bremerhaven, Germany. <sup>9</sup>Helmholtz Institute for Functional Marine Biodiversity, Oldenburg, Germany. <sup>10</sup>These authors contributed equally: Natalie J. Wu-Woods, Jacob T. Barlow. ✉ e-mail: [rustem.admin@caltech.edu](mailto:rustem.admin@caltech.edu)

profiles of microbial communities. More recently, shotgun metagenomics—sequencing of the entire DNA content of a sample—has become more common in human microbial ecology due to its ability to provide in-depth, genome-resolved characterizations of microbial populations<sup>12,13,15,16</sup>. Further, by resolving subpopulation structures within a single taxon, shotgun metagenomics can enable additional insights, such as how physiological host gradients can induce evolutionary pressures on the microbiome<sup>17,18</sup>. Genome-resolved characterization is also needed to study which microbial genes are under selection pressure under different host environments.

The molecular details of how tissue-associated microbes interact with the host environment remains poorly understood because the field lacks the appropriate tools to go beyond taxonomic profiling and investigate microbial pathways and genes directly from intestinal biopsies. Two common methods of full-genome characterization include culturing microbial isolates or reconstructing metagenome-assembled genomes (MAG) directly from mixed microbial samples. Culture-dependent methods have their role; however, culture-independent methods are attractive for characterizing microbes from their native context, as well as those that cannot be easily isolated. MAGs are created via a computational approach<sup>15</sup> in which sequencing reads are assembled into continuous sequences and then grouped into separate bins to reconstruct complete genomes without culturing<sup>19</sup>.

Shotgun metagenomic sequencing analyses of complex host-associated microbiomes have been challenged by the high ratio of host to microbial nucleic acids<sup>20</sup>. In humans, 85–95% of reads in a saliva sample<sup>21</sup> are host and more than 99.99% of reads in an intestinal biopsy are host. These enormous ratios of host to microbial DNA (1:10,000 in a human intestinal biopsy) are particularly challenging for shotgun metagenomic sequencing studies because most reads align to the host genome. Such tissue samples sequenced directly using current protocols and sequencing depths do not produce sufficient microbial reads to construct MAGs.

To prevent most shotgun-sequencing reads assigning to host, a wide variety of host removal (that is, host-depletion) methods have been developed<sup>21–26</sup>. Published and commercial protocols have enabled both long-read sequencing and bacterial MAG construction from mammalian derived liquid samples. Although some protocols have been validated for use on solid-tissue sample types, and others may have potential for success in these samples, none have been shown to be sufficiently effective to enable bacterial MAG construction from solid mammalian tissues. Additionally, many host-depletion methods are not feasible to perform in the clinic due to extensive processing times and complex protocols.

Here, we developed and optimized a microbial-enrichment method (MEM) to remove host nucleic acids from complex samples while not substantially perturbing the microbial community composition. We demonstrate the performance of MEM in laboratory and clinical settings, and with a range of sample types, including saliva, feces, intestinal scrapings and intestinal biopsies. We also demonstrate the ability of MEM followed by shotgun metagenomic sequencing to detect both high- and low-abundance microbial taxa, pathways and genes from human intestinal biopsies along the GI tract. Moreover, we show the use of MEM to enable MAG construction directly from human intestinal biopsies to identify and differentiate subpopulations and subpopulation variants.

## Results

### MEM development

We developed a MEM that incorporates a selective-lysis protocol using mechanical stress (bead beating) by leveraging the size differences between host and bacterial cells (Fig. 1a,b). Beads typically used for microbial lysis are 0.1–0.5 mm but we chose larger beads (1.4 mm) to create high mechanical shear stress on the larger host cells while leaving

small bacterial cells intact<sup>27</sup>. Next, Benzonase is added to degrade accessible extracellular nucleic acids, including nucleic acids from dead lysed microbes. Proteinase K further lyses host cells and degrades host histones for DNA release. We also tested and optimized other factors affecting the performance of MEM, including enzymatic nucleic acids removal, bead beating and incubation time, to keep the entire protocol time under 20 min, with gentle processing conditions to prevent microbe lysis. (Methods and Supplementary Fig. 1).

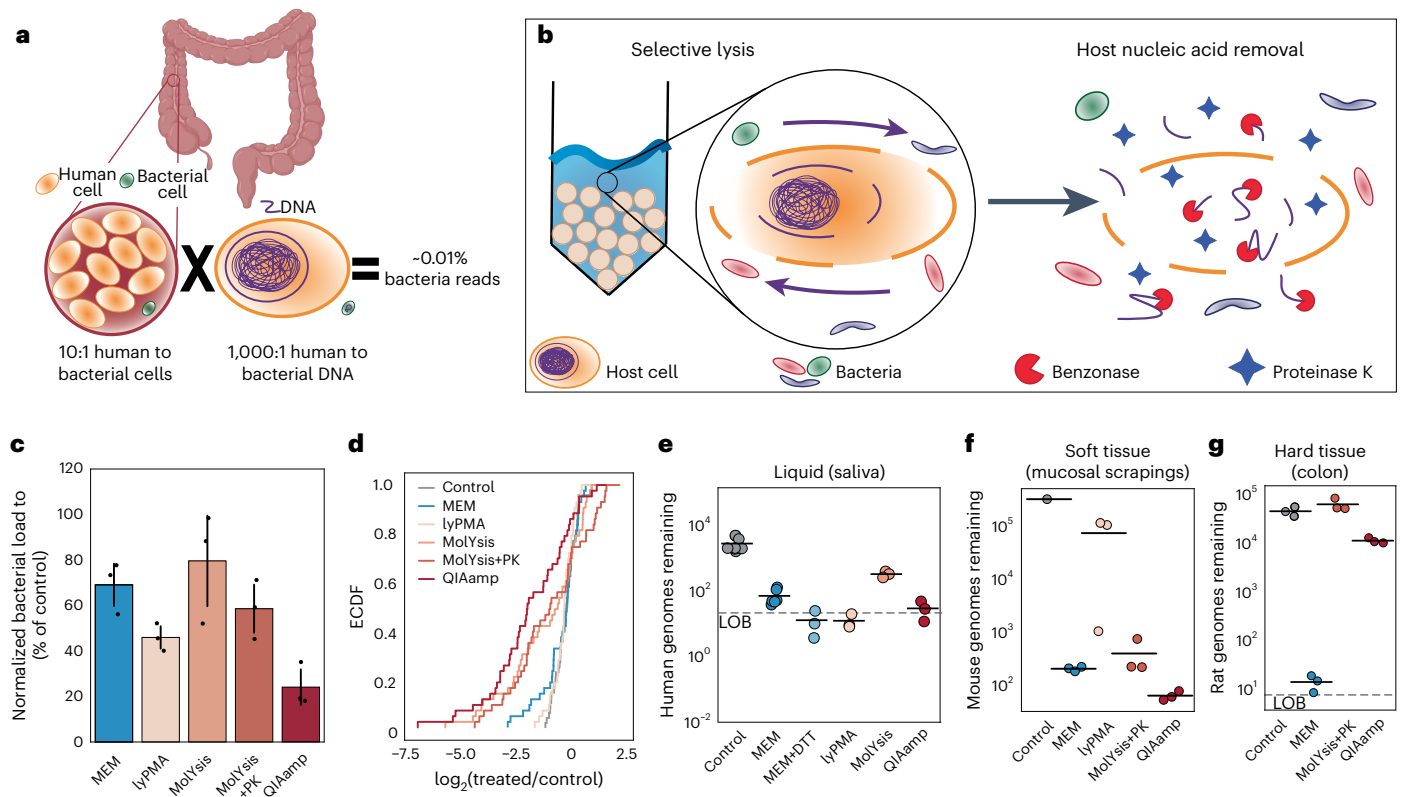
To compare host depletion by MEM with existing methods, we selected three published methods that use different cell-lysis approaches: MoLYsis, QIAamp and lyPMA. All host-depletion methods include two main steps: selective lysis followed by nucleic-acid removal. QIAamp lyses cells lacking a cell wall through a weak detergent, saponin<sup>28</sup>. MoLYsis selectively lyses the more fragile mammalian cells through exposure to a weak concentration of guanidinium<sup>29</sup>. lyPMA lyses mammalian cells through osmotic lysis<sup>21</sup> and uses photochemistry to render DNA accessible to propidium monoazide (PMA) nonamplifiable.

### MEM minimizes loss of bacteria during sample processing

To quantify how MEM affects microbial community composition and relative abundances of individual taxa, we first used frozen mouse fecal samples. We chose fecal samples instead of a contrived community to characterize microbial effects on a range of unique taxa and on a continuum of abundances. This is because mouse fecal samples do not typically require host depletion because they have low levels of host contamination (more than 90% of the DNA biomass originates from nonhost cells). The high biomass of microbial cells makes feces ideal for characterizing the impact of different host-depletion methods on the microbial community composition. Additionally, contrived communities still require an extracted control due to variation in extraction kit efficiency<sup>30,31</sup>.

Although we were unable to extract all DNA molecules in the samples, all samples were processed with the same extraction kit following host depletion to standardize extraction kit and/or lysis efficiency. On homogenized stool samples, we observed similar losses in microbial recovery across all five host-depletion protocols compared to a control, untreated sample (Fig. 1c). MEM induced on average 31% (standard deviation (s.d.) 11%) bacterial loss, which falls within the expected fraction of 10–50% dead microbial cells in stool<sup>32</sup>. To characterize how MEM and the other host-depletion methods affect the microbiome at a taxonomic level, we next performed quantitative 16S rRNA gene sequencing<sup>33</sup> on the mouse fecal samples ( $n = 3$ ). By comparing paired host-depleted and control samples, we found that lyPMA and QIAamp induced the largest total bacterial losses whereas MoLYsis and QIAamp induced the least uniform bacterial losses, with some taxa dropping more than 100-fold (Fig. 1c,d). Previous literature suggests QIAamp's saponin concentration can be lowered to limit some of these bacterial losses<sup>25</sup>. We confirmed that MEM induced minimal losses in the microbial community; more than 90% of genera showed no significant difference in relative abundance between MEM and control samples (paired  $t$ -test, two-sided,  $P = 0.05$ ). Additionally, all taxa that were consistently detected in the control samples were also detected in the MEM-treated samples, whereas MoLYsis and QIAamp resulted in some taxa drop out (Supplementary Table 1). Because MEM selectively lyses host cells based on cell size, this approach appears to introduce lower bacteria bias compared with chemical lysis alternatives (MoLYsis and QIAamp) where degree of lysis may differ based on bacterial cell wall and/or membrane structures (Supplementary Fig. 2).

To determine how effectively MEM and the other host-depletion methods removed host material, we next quantified the amount of host DNA remaining after each host-depletion method on three additional sample types: liquid, soft-tissue and hard-tissue samples (Fig. 1e–g, discussed in detail below), in which the host DNA made up as much as 99.9% of the total biomass. In saliva, all methodologies enabled some



**Fig. 1 | Comparison of the performance of the MEM with published host-depletion methods.** **a**, Estimated percentage of bacterial reads obtained when human intestinal biopsies are sequenced without processing. **b**, Schematic demonstrating the two-step selective-lysis and nucleic-acid removal techniques used in MEM. **c**, Bacterial loads from mouse stool samples treated with five different host-depletion methods. Loads are normalized to the control (no host depletion) stool samples ( $n = 3$ ; error bars are 95% confidence interval centered on the mean). **d**, Empirical cumulative distribution function (ECDF) of 16S rRNA gene amplicon sequencing results from mouse stool samples normalized to the control stool samples ( $n = 3$ ). Curves shifted to the left of the control indicate a greater percentage of taxa with lower abundance than the control samples following host depletion. **e–g**, Remaining host DNA was quantified through ddPCR of a single-copy host specific primer (Methods). Reported genomes

remaining refers to the abundance of this single-copy gene present in 1  $\mu$ l of elution. **e**, Remaining human genomes in fresh human saliva were quantified after treatment with each host-depletion method and in untreated controls ( $n = 3$  biological replicates for lyPMA, MolYsis and QIAamp;  $n = 4$  biological replicates for Control;  $n = 4$  biological replicates for MEM and  $n = 3$  technical replicates for MEM + DTT). **f**, Host-depletion methods were tested on mouse intestinal mucosal scrapings as a representative of soft tissue and remaining mouse genomes were quantified ( $n = 3$  biologic replicates for host depletion methods and  $n = 1$  for control from one mouse). **g**, Host-depletion methods were tested on rat colonic sections as a representative hard tissue (including connective tissue, muscle and mucosa) and remaining rat genomes were quantified ( $n = 3$ ; biologic replicates from one rat).

host removal. Following MEM treatment, over 40-fold depletion of host was achieved (Fig. 1e). The addition of dithiothreitol (DTT) pretreatment, which was added due to the high mucin content of saliva, slightly increased host removal by MEM in some participants (Fig. 1e and Supplementary Fig. 3). lyPMA appeared highly effective at host removal, but was difficult to use predictably as the stoichiometric nature of the method can result in large microbial losses when host levels are lower than expected (Supplementary Fig. 3). Additionally, MolYsis showed increased bacterial recovery, likely due to the additional mutanolysis step<sup>22,26</sup>.

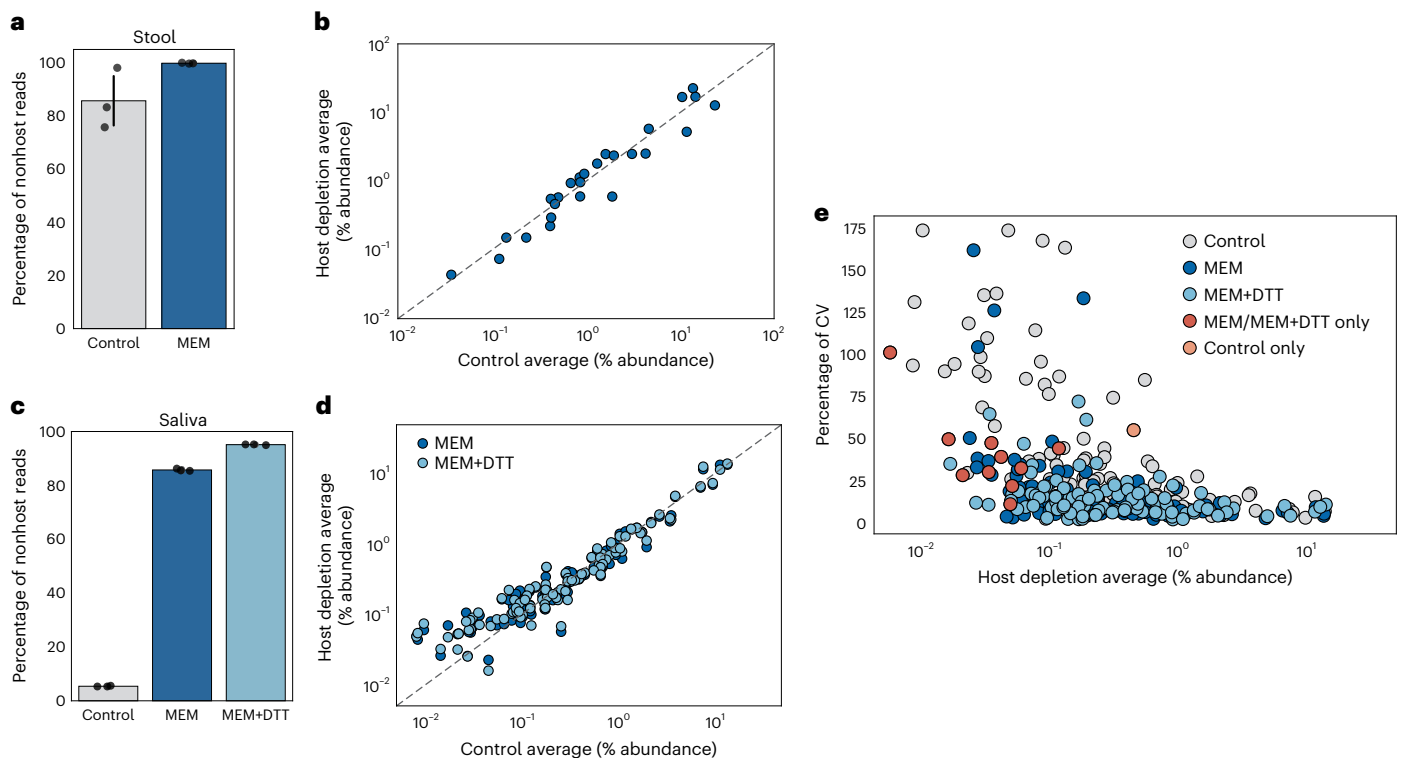
We next examined host depletion on whole tissue samples, beginning with mouse intestinal scrapings that isolate the epithelial layer with mucosa-associated bacteria (Supplementary Fig. 4). Mucosal scraping samples were efficiently host-depleted by MEM and some of the published methods (Fig. 1f). MEM, MolYsis and QIAamp all showed around 1,000-fold depletion of host with QIAamp showing slightly greater host removal (MEM had an average 1,600-fold depletion, s.d. 170). lyPMA performed poorly on the soft-tissue sample because this method relies on ultraviolet-activated crosslinking making it incompatible with opaque sample types.

We next tested host-depletion methods on hard-tissue samples using rat colonic sections because they are anatomically similar to a human intestinal biopsy. We excluded lyPMA from the hard-tissue experiment due to its poor performance on soft tissue (Fig. 1f). MEM was the only method that worked on the solid-tissue sample type (Fig. 1g). MEM treatment resulted in almost complete removal of the host DNA (3,600-fold removal, s.d. 1,500), whereas MolYsis and QIAamp host DNA levels after treatment were similar to the control.

These experiments demonstrated that MEM is a solid-tissue host-depletion method that can remove host DNA more than 1,000-fold while introducing minimal losses in the relative abundances of the microbial fraction. In our experiments, more than 90% of genera showed no significant difference (paired  $t$ -test, two-sided,  $P = 0.05$ ) in relative abundance between MEM-treated and control samples based on 16S rRNA gene profiling.

### Shotgun sequencing of MEM-treated saliva and stool

We next investigated the use of MEM for shotgun metagenomics. Owing to biomass limitations, accurate characterization of microbial communities through shotgun sequencing of control (not host-depleted)



**Fig. 2 | Microbial enrichment of stool and saliva after host depletion by MEM as confirmed by shotgun sequencing.** **a**, The percentages of nonhost reads in control and MEM-treated mouse stool samples were calculated bioinformatically through alignment to a mouse reference genome ( $n = 3$ ; error bars are 95% confidence interval centered on the mean). **b**, Species-level taxon relative abundances were plotted for control and MEM-treated mouse stool and overlaid on a dashed line showing a 1:1 correlation. **c**, Shotgun sequencing was performed on control and MEM-treated fresh human saliva. The percentages of nonhost reads were calculated bioinformatically through alignment to a human reference genome. One saliva sample was evenly split nine ways for this comparison ( $n = 3$ ). **d**, Species-level taxon relative abundances were plotted

for control and MEM-treated fresh human saliva and overlaid on a dashed line showing a 1:1 correlation. An additional DTT pretreatment was performed before MEM treatment for a subset of MEM-treated samples (MEM + DTT) ( $n = 3$ ). **e**, Coefficient of variation (CV) was plotted against relative species abundance and colored based on treatment types in which the taxa were detected. Each point represents a species; gray, dark-blue and light-blue points indicate taxa that were present in all three treatments (control, MEM and MEM + DTT). MEM/MEM + DTT only (red points) indicate the ten taxa found only in the MEM-treated samples. Control only (orange point) indicates the single taxon that was found only in the control samples, which was identified as *Haemophilus*.

samples is not feasible on intestinal biopsies. Thus, we first used saliva and stool samples to investigate potential biases associated with MEM treatment within the microbial fraction. We first confirmed that MEM treatment enabled reliable reduction in host reads across mouse stool and human saliva samples (Fig. 2a,c). DTT pretreatment in saliva improved host removal roughly tenfold (Fig. 2c).

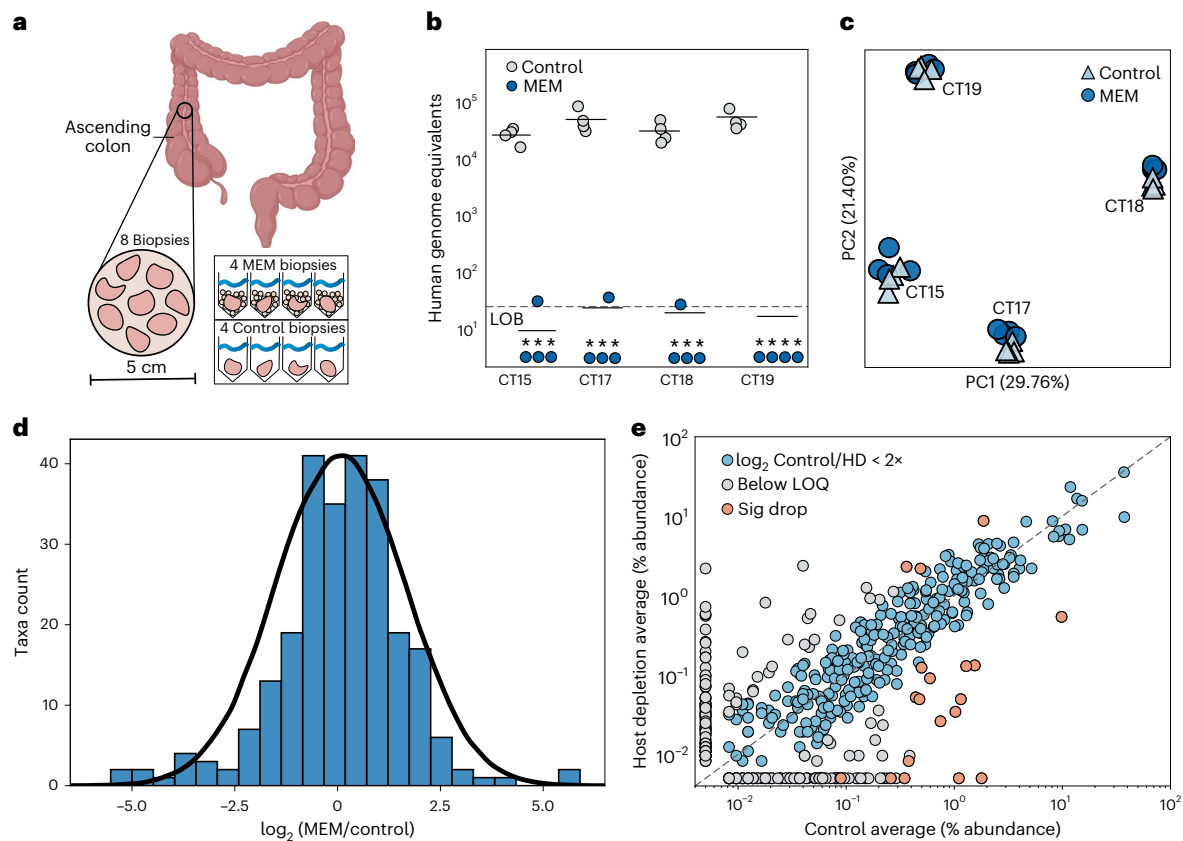
Next, we compared the results of shotgun sequencing between the control and MEM-treated samples. There was a high correlation between the relative abundances of bacterial taxa in control and MEM-treated samples for both stool and saliva (Pearson coefficient of determination,  $R^2 = 0.93$  for stool and  $R^2 = 0.90$  for both MEM and MEM + DTT in saliva; with  $R^2 = 0.93$  for taxa above 0.1% relative abundance). For stool and saliva, a high correlation between the relative abundance of species in control versus MEM-treated samples showed that MEM did not substantially alter microbiome composition (Fig. 2b,d and Supplementary Tables 2 and 3). For saliva, the correlation was less pronounced for low-abundance taxa with enrichment of specific species in MEM-treated samples and was investigated more quantitatively by comparing the coefficient of variation across samples for low-abundance species (Fig. 2e). MEM-treated saliva samples had lower coefficient of variation (50% coefficient of variation, 95% confidence interval), indicating better replicability compared with untreated controls. Additionally, MEM improved quantification of low-abundance species (0.05–0.5% relative abundance), enabling detection of an additional ten species that were undetected in the

control. We further confirmed these taxa were not introduced during MEM processing (Supplementary Fig. 5). These shotgun-sequencing experiments with mouse fecal and human saliva samples demonstrated that MEM treatment introduced minimal microbial biases (more than 98% of microbial species experienced less than a fourfold loss in relative abundance) while detecting additional microbial taxa at equivalent sequencing depths.

### MEM feasibility on human intestinal mucosal biopsies

To determine how MEM performs on human intestinal biopsies, we recruited healthy participants undergoing routine colon cancer screenings via colonoscopy. From each of four participants, we obtained eight mucosal biopsies; four biopsies from each participant were assigned to the MEM-treatment group and four were untreated controls (Fig. 3a). Owing to concerns regarding contaminant DNA in samples with low bacterial loads<sup>34–37</sup>, we also characterized the background bacterial signal associated with MEM and our processing methods through quantitative 16S rRNA gene sequencing of MEM processing blanks (Methods and Supplementary Tables 4 and 5). MEM removed host DNA more than 2,000-fold across all 16 biopsies, with most biopsies having host levels comparable to a processing blank after MEM treatment (Fig. 3b).

To determine how MEM affects the human intestinal microbiome at a community level, we first performed quantitative 16S rRNA gene sequencing<sup>33</sup>. Roughly 93% of genera remained in MEM-treated and control biopsies after computationally removing taxa found at higher



**Fig. 3 | Analysis of microbial enrichment in paired human intestinal biopsies processed with and without MEM.** **a**, Illustration of sample collection. **b**, Host DNA was quantified for each biopsy using ddPCR of a single-copy host primer. Human genomes remaining refers to the abundance of this single-copy gene present in 1  $\mu$ l of elution (\* indicate measurement was below limit of blank, LoB). **c**, Biopsies were characterized with 16S rRNA gene sequencing and principal component analysis (PCA) on microbial genus-level relative abundances were performed to visualize microbial population variation. **d**, The  $\log_2$ -fold

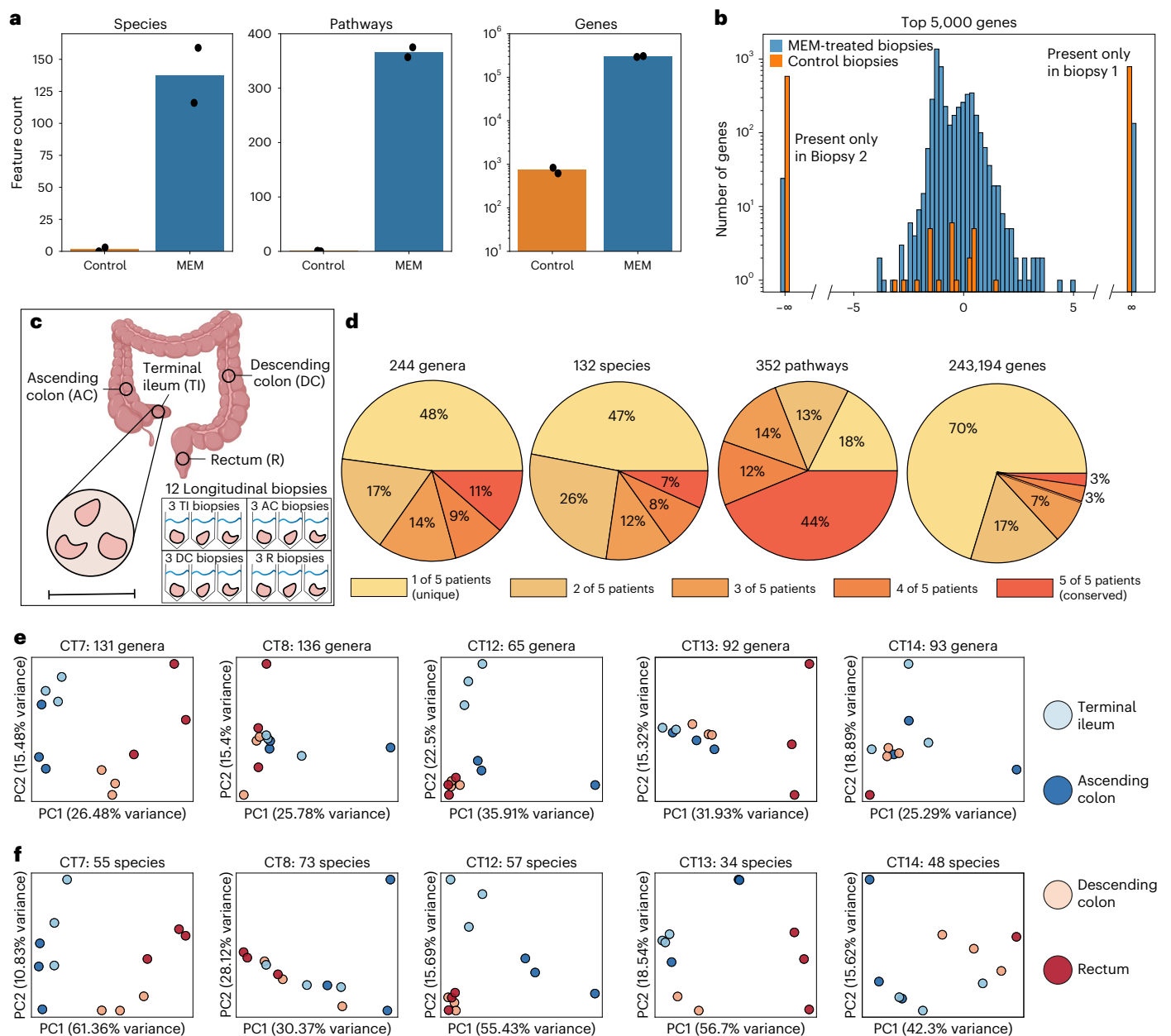
differences on microbial genus-level relative abundances between control and MEM-treated biopsies were plotted with a standard normal distribution overlaid in black. **e**, Microbial genus-level relative abundance measured in control versus MEM-treated biopsies were plotted and overlaid on a dashed line showing a 1:1 correlation. Highlighted in gray are taxa that were below the assay limit of quantification (LOQ). Highlighted in orange are taxa with greater than fourfold changes between control and MEM biopsies.

absolute abundances in the blanks, giving us confidence most detected taxa were not background contaminants. To further confirm that MEM did not introduce additional contamination, we found strong agreement between taxon abundances in MEM and control biopsies (Supplementary Table 6). Principal component analysis of sequencing results showed that any differences in microbial relative abundances introduced by MEM were less than the differences observed between participants (Fig. 3c). Analysis of sequencing results revealed minimal changes in relative abundances of most taxa after MEM treatment, with roughly 88% of taxa having no significant differences in relative abundances from the controls (Mann–Whitney  $U$ -test, two-sided  $P = 0.05$ ). For taxa present at greater than 1% relative abundance, more than 95% of taxa had no significant differences between MEM and control samples. The  $\log_2$ -fold difference in taxa between control and MEM-treated samples approximated a normal distribution (Kolmogorov–Smirnov test against normal distribution, statistic of 0.074,  $P = 0.11$ ) (Fig. 3d). Furthermore, there was a linear correlation in relative taxon abundances between the control and MEM-treated samples (Fig. 3e). MEM enables over 1,000-fold host removal while introducing minimal biases in microbial relative abundances when used in a clinical setting on human intestinal biopsies.

### MEM enables study of microbial species, pathways and genes

To investigate whether MEM enables detection and characterization of additional microbial species, pathways and genes from

human intestinal biopsies, we shotgun sequenced paired control and MEM-treated biopsies from CT18 at a depth of above 100 million reads ( $n = 2$  for each condition) (Fig. 3a). We observed a roughly 100-fold increase in the number of organisms, a 700-fold increase in the number of pathways and more than a 400-fold increase in the number of genes detected in MEM-treated samples compared with the control samples (Fig. 4a). When comparing only completed pathways, defined as above 90% complete, no complete pathways were detected in either of the control biopsies. An average of 1.5 ( $\pm 1.5$ ) species and 728 ( $\pm 107$ ) genes were detected in the control biopsies, whereas an average of 137.5 ( $\pm 21.5$ ) species and 300,641 ( $\pm 6,922$ ) genes were detected in the MEM-treated biopsies. MEM treatment enabled shotgun-sequencing classification of microbes down to a relative abundance of 0.005%, whereas in control biopsies a minimum relative abundance of 10% was required to detect microbes at similar sequencing depths. MEM-treated biopsies could detect genes down to a relative abundance of  $10^{-10}$ , whereas in control biopsies genes could only be detected when present at a minimum relative abundance of  $10^{-5}$ . We further found that MEM treatment improved reproducibility of detecting the most abundant genes. We compared the relative abundance of the top 5,000 most abundant genes between two MEM-treated and two control biopsies (Fig. 4b). In the control biopsies, a high percentage of the genes (98%) were detected in only one sample, whereas for the MEM-treated samples only 3% of the genes were detected in one biological replicate but not the other.



**Fig. 4 | Shotgun sequencing of MEM-treated human intestinal biopsies.**

**a**, Four biopsies from participant CT18 (two MEM-treated and two control) were shotgun sequenced and the number of microbial species, pathways and genes identified in each sample were plotted ( $n = 2$ ). **b**, For the top 5,000 abundant genes, the log<sub>2</sub>-fold-change in relative abundances between the two MEM-treated biopsies and the two control biopsies were plotted. **c**, Illustration of sampling collection. Scale bar, 5 cm. **d**, The number of features for genera,

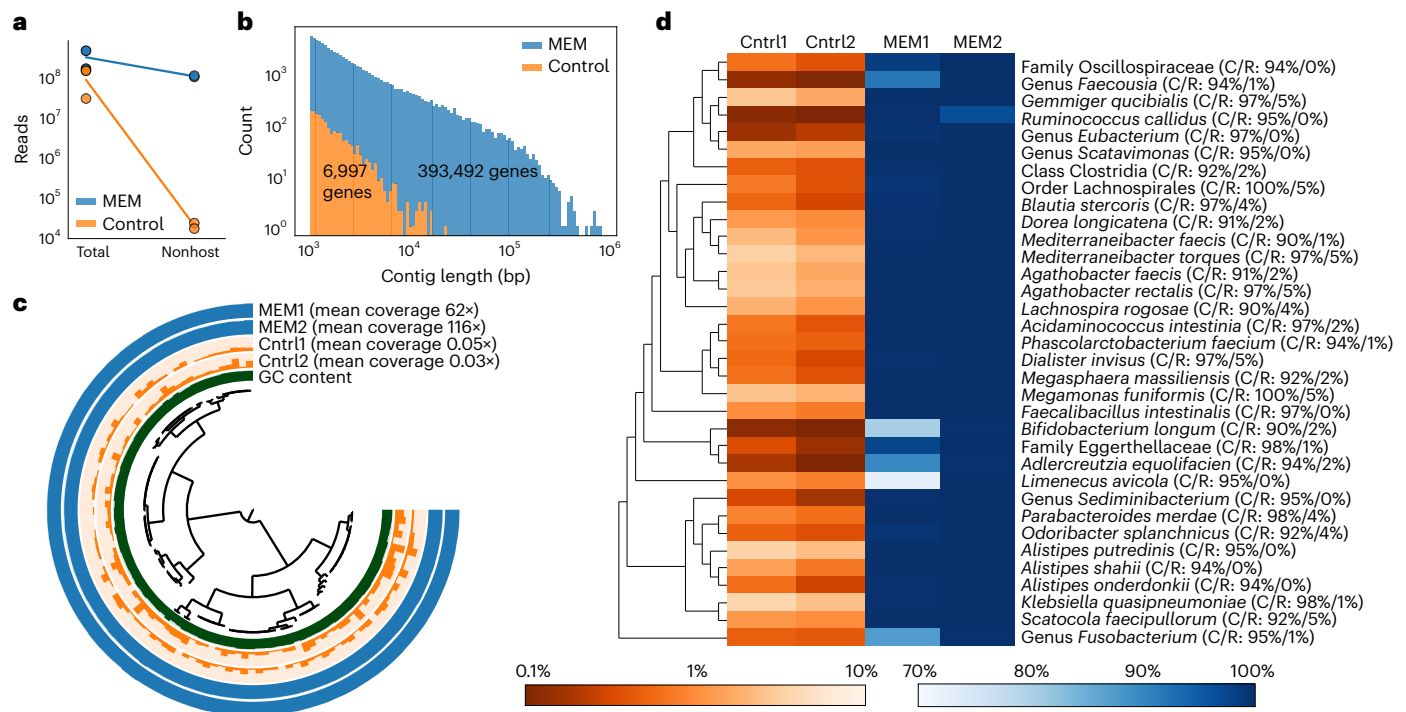
species, pathways and genes were grouped based on whether they were present in at least one biopsy sample from only one participant, two participants, three participants, four participants or from all five participants. **e, f**, Principal component analysis (PCA) was carried out on all 60 longitudinal samples grouped by participant. **e**, Principal component analysis on relative abundance of 16S rRNA gene sequencing genera assignments. **f**, Principal component analysis on relative abundance of shotgun-sequencing species assignments.

Next, we tested whether MEM would enable characterization of microbial variation (at the taxon, pathway and gene levels) cross sectionally across individuals and longitudinally across the GI tract of a single individual. Five healthy participants undergoing colonoscopy were sampled in four regions across the GI tract: terminal ileum, ascending colon, descending colon and rectum. From each location, three biopsies were obtained resulting in a total of 12 biopsies per participant (Fig. 4c). All biopsies were processed with MEM and the microbial profiles were characterized via 16S rRNA gene sequencing and shotgun sequencing at an average read depth of 25 million, producing an average of 2 million nonhost reads (Fig. 4d–f and Extended Data Figs. 1 and 2). About half (91 of 187) of the microbial species identified

were unique to an individual (Fig. 4d and Supplementary Table 7). These unique species ranged in relative abundance from 10 to 0.01% (Supplementary Fig. 6). As was observed previously, pathways appeared more conserved across participants compared with taxonomy (genera and species)<sup>38,39</sup>.

#### Variation in mucosal microbes longitudinally in the GI tract

Whether mucosa-associated microbes vary along the GI tract has been challenging to determine due to the low number of microbial reads that could be recovered from mucosal biopsies<sup>12</sup>. We first tested whether microbial variation between GI sites is present at the genus-level. For each participant sample, we used quantitative<sup>33</sup> 16S rRNA gene



**Fig. 5 | MAG construction with MEM-treated human intestinal biopsies performed from shotgun metagenomic sequencing. a**, Two control and two MEM-treated biopsies from the same participant (CT18) and intestinal region (ascending colon) were shotgun sequenced. The number of nonhost reads was determined after alignment to a human reference genome. **b**, Contigs were constructed from coassembly of the two samples from each condition and the distribution of contig lengths was plotted. The number of prokaryotic genes identified in these contigs is shown. **c**, MAG of *Alistipes putredinis* was constructed from coassembled MEM biopsies. Bar heights represent mean coverage and are scaled independently for each sample. **d**, From coassembly

of MEM biopsies, 34 high-quality MAGs (more than 90% complete, less than 5% redundant) were constructed de novo. The heatmap shows the percentage of each genome that is covered by at least once by the sample (that is, detection or breath of coverage), with a maximum of 3.7% in control samples and 99.999% in MEM samples. The average detection for MEM1, MEM2, Cntrl1 and Cntrl2 were 97.3% (s.d. 6.4%), 99.8% (s.d. 0.7%), 1.2% (s.d. 1.1%), and 0.8% (s.d. 0.7%), respectively, across all MAGs. Taxonomy was assigned for each MAG and listed to the right along with completion and redundancy (C/R). The phylogenetic tree to the left of the heatmap highlights taxonomic grouping of each MAG.

sequencing to quantify genus-level microbial changes longitudinally along the GI tract. Microbial taxa from the proximal colon (terminal ileum and ascending colon) and taxa from the distal colon (descending colon and rectum) showed some clustering by location in most participants (Fig. 4e). Each participant sample was shotgun sequenced to test whether the observed variation in taxa along the GI tract extended to the species, pathway or gene levels. Clustering between the terminal ileum and ascending colon versus the descending colon and rectum was seen in some participants across species, namely in participants CT7, CT12, CT13 and CT14 (Fig. 4f). There appeared to be minimal clustering between the terminal ileum and ascending colon versus the descending colon and rectum at the pathway and gene-level (Extended Data Fig. 3). Additionally, there was high variation within regions for some individuals, which may be attributed to read depth limitations. For example, for one descending colon sample from CT13 no microbial marker genes were identified due to the minimal number of nonhost reads (Fig. 4f).

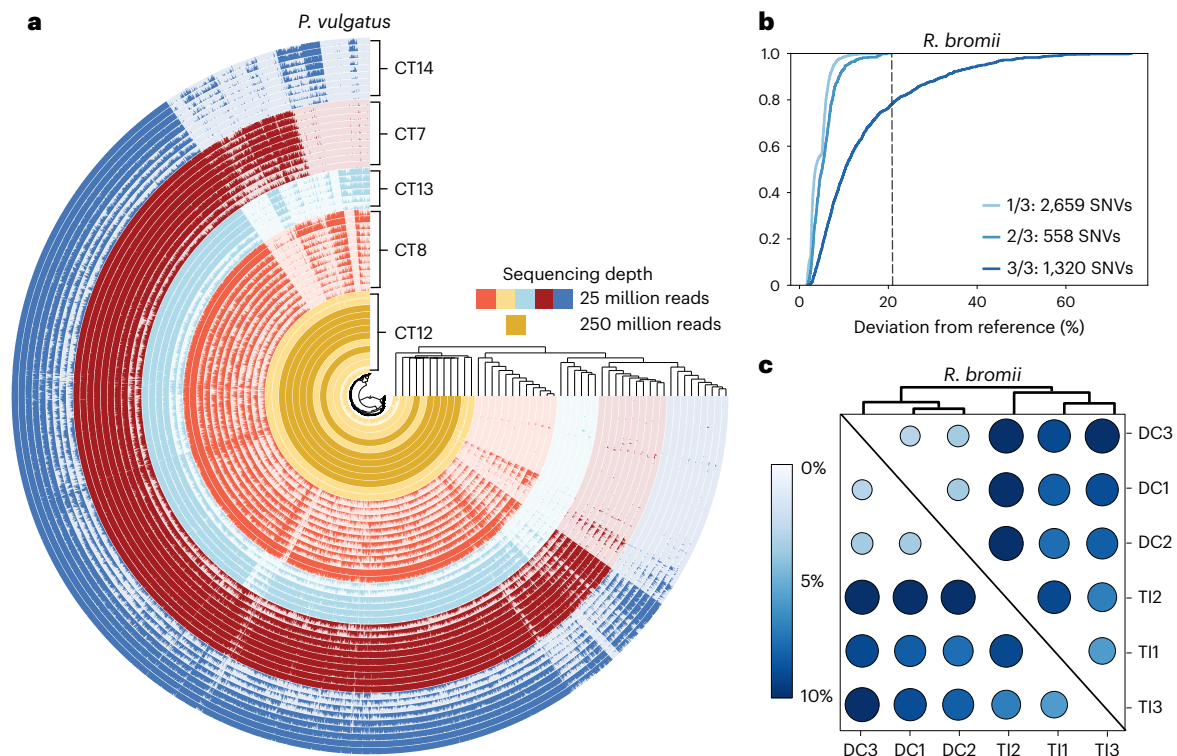
Shotgun sequencing of MEM-treated human intestinal biopsies enabled characterization of high- and low-abundance microbial species, pathways and genes. This characterization documented longitudinal shifts in the mucosal microbiome along the lower human GI tract. To investigate whether a single microbial strain varies along the GI tract, we next attempted to assemble microbial genomes from MEM-treated intestinal biopsies.

### MEM enables MAG of intestinal microbes from human biopsies

To determine whether MAGs could be constructed after processing with MEM, we selected two control and two MEM-treated biopsies with similar bacterial loads from participant CT18 (Figs. 3a and 4a). Samples

were shotgun sequenced and processed for genome reconstruction as previously described<sup>15</sup>. We sequenced both control and MEM-treated biopsies to measure the additional information MEM treatment can help yield at equivalent sequencing depths. After processing, host reads were removed bioinformatically and roughly 10% of reads were identified as nonhost in MEM-treated samples whereas roughly 0.01% of reads were identified as nonhost in the untreated controls (Fig. 5a).

We first tried to reconstruct MAGs from the control samples, however, the assembly of the short reads from nonhost-depleted samples and our subsequent attempts to bin the resulting contigs into MAGs were unsuccessful because these assemblies suffered from remarkably short contigs (Fig. 5b). Coassembly was then performed on MEM-treated samples and resulted in substantially more and longer contigs compared with the control samples, with contig lengths of up to 833 kbp (Fig. 5b and Supplementary Table 8). Automatic binning and manual refinement steps resulted in a total of 34 high-quality bacterial MAGs (more than 90% complete and less than 5% redundant) and more than 70 medium-quality MAGs (more than 50% complete and less than 10% redundant), demonstrating how MEM treatment of human intestinal biopsies makes it possible to reconstruct MAGs from these samples. For the 34 high-quality MAGs, we computed detection, which reports the proportion of nucleotides in a given reference sequence that are covered by at least one short read in a given metagenome. Thus, detection is an extremely effective way to be able to discuss the presence of a given population in a given sample, independent of read coverage, and by avoiding false positives due to nonspecific read recruitment. To confirm that the MAGs reconstructed from MEM-treated samples were accurate representations of the untreated biopsies, we assessed



**Fig. 6 | Interindividual and intraindividual bacterial biodiversity present along GI tract. a**, Gene-level analysis was carried out on *P. vulgatus* for all five participants. Samples were grouped by gene detection, defined as percentage of each gene with at least 1× coverage, and showed strong participant-dependent grouping but lacked grouping by GI location. **b**, Empirical cumulative distribution function of the occurrence of SNVs in a MAG of *R. bromii* and the deviation of these SNVs from the reference across three technical replicates. 1/3, 2/3 and 3/3 indicate the number of technical replicates that had an SNV at that

location followed by the total number of SNVs in each of these categories. A black dashed line is drawn at 21% deviation from reference; above this value, all observed SNVs were present in all three technical replicates. **c**, Nucleotide-level analysis was carried out on MAGs with a mean coverage above 50× across all samples. Shown here is the fixation index from SNVs analyzed within the coding region of *R. bromii* with a minimum deviation from reference set at 21%. Samples were clustered based on fixation index and strong region-dependent groupings can be seen. DC, descending colon; TI, terminal ileum.

the uniformity of coverage of the control reads when mapped back onto the MAGs (Fig. 5c). To perform this analysis, we chose a MAG resolved to *Alistipes putredinis*, a known gut microbe that had the highest detection in the control samples. Control samples showed an even distribution of reads among the 29 contigs present, indicating that this MAG was also present in the control samples, but sequencing depth limitations prevented the reconstruction of a genome. Overall, we observed a higher detection of all 34 high-quality MAGs in MEM-treated samples compared to control samples (Fig. 5d).

To quantify whether reads from control samples mapped back onto all 34 high-quality MAGs, detection was plotted for each MAG (Fig. 5d). Next, to assess whether any of these MAGs were contaminated, we performed taxonomic classification on each genome<sup>40,41</sup>. With a threshold of 95% average nucleotide identity (ANI), 33 MAGs were successfully classified. We compared the size of each classified MAG with the matching reference genomes in the Genome Taxonomy Database (GTDB) and found high agreement with current microbial databases ( $R^2 = 0.78$ ,  $P = 4.32 \times 10^{-12}$ ) indicating that the MAGs constructed from MEM-treated samples were not artifacts (Supplementary Tables 9 and 10). One *Fusobacterium* MAG matched closely with a published fecal-derived MAG at 86.85% ANI, but GTDB was unable to assign species-level taxonomy (Supplementary Fig. 7). Because all MAGs were constructed in the same manner and with similar quality metrics, it is likely that this *Fusobacterium* MAG is a uncharacterized taxon rather than contamination. We also wanted to quantify the range of microbial diversity we could capture with MAGs. These 34 MAGs spanned six bacterial phyla (Fig. 5d) and an archaeon

(*Methanobrevibacter smithii*) MAG was constructed from participant CT12, demonstrating MEM-treated biopsies enabled genome reconstruction of archaea and a wide variety of bacteria (Extended Data Fig. 4). Using MEM, high-quality microbial MAGs were reconstructed from microbes from human intestinal biopsies at relative abundances down to 1%.

#### MEM identifies distinct microbial strains across individuals

After establishing the feasibility of MAG construction directly from human intestinal biopsies, we next investigated how microbial genomes may vary across individuals and within individuals. To determine whether MEM enables differentiation of population-level microbial differences across individuals within a single taxon, a total of six biopsies from participant CT12 were resequenced to a sequencing depth of roughly 250 million reads. Assembly and binning were performed on each of the six biopsies individually and MAGs were dereplicated across samples. A MAG of *Phocaeicola vulgatus*, the most prevalent and abundant species found in all participants, from participant CT12 was constructed and annotated (Fig. 6a). Reads from all 60 intestinal biopsies taken from all five participants were then mapped onto this MAG to identify which genes were absent from the other participants' biopsies.

A gene-resolved analysis of naturally occurring *P. vulgatus* populations through metagenomic read recruitment was performed, as described previously<sup>42</sup>, and revealed a large core genome and differentially occurring genes across individuals (Fig. 6a). Genes from biopsies taken across GI tract regions within participant CT12 appeared conserved (CT12 samples had an average gene detection of more than 96%).



Some genes with high detection were only found in one or two participants (either CT12 only or CT12 and one other participant), which we defined as unique genes. To assess whether these genes were functionally distinct, genes were annotated with the Clusters of Orthologous Groups (COG) database to identify orthologous genes. Of the 287 genes unique to CT12, 100 of these were annotated by COG and corresponded to a wide range of functions (Supplementary Fig. 8). Of the gene clusters unique to two participants (that is, CT12 and one other individual), about 30% were annotated (Supplementary Fig. 8). MEM treatment enables insights into functionally distinct microbial populations of the same taxon that occupy the same geographical location in the gut across individuals with similar health status.

### SNVs detectable across GI tract regions within an individual

Finally, we investigated whether MEM treatment could enable studies of microbial population genetics in low-biomass samples through single-nucleotide variants (SNVs) as a result of the increased depth of coverage. For this analysis, we analyzed MAGs from participant CT12 as reference genomes and mapped the paired-end reads from the terminal ileum and descending colon from CT12 onto these assembled genomes. Six MAGs had a mean coverage above 50 $\times$  across all six samples (three terminal ileum and three descending colon) and were selected for subsequent SNV analysis (Supplementary Fig. 9). SNV profiles were generated from the paired-end reads of each sample by comparing them with the reference sequence (MAG). We investigated whether PCR errors may be responsible for some of the SNVs observed in our data by preparing libraries for an additional three technical replicates from a single terminal ileum biopsy (Fig. 6b), with the expectation that differences in the SNV profiles of the technical replicates should be minimal. By looking at nucleotide variations occurring in one, two or all three replicates, we observed that a minimum deviation from the reference nucleotide of 21% for *Ruminococcus bromii* (Fig. 6b) allowed for the selection of SNVs only and minimized the impact of PCR errors in the population structure analysis. Analyses of these data using fixation index showed that some taxa, such as *R. bromii* (Fig. 6c) and *Gemmiger formicilis* (Supplementary Fig. 9), were composed of subpopulations that were distinct between the upper and lower intestinal tract. To assess whether these SNVs were functionally important, codon-level and translated (amino acid) analyses of SNVs in *R. bromii* were performed and similar clustering of biopsies by location was detected (Extended Data Fig. 5). The recovery of SNVs afforded by the deeper sequencing and increased coverage of MAGs from biopsy samples allowed us to detect the presence of subpopulation structures for some individual taxa along the lower GI tract of a single individual.

### Discussion

MEM is a method for use on mammalian host-rich sample types that enables metagenome shotgun sequencing and analysis of microbes present in these samples. MEM enables more than 1,000-fold removal of host DNA from solid mammalian tissue while minimally perturbing the microbial community composition. MEM is simple and fast, with processing times less than 30 min, facilitating integration into a laboratory or clinical workflow without in-person training. MEM is highly compatible with downstream shotgun sequencing of microbial DNA, leading to the detection of more than 400-fold more species and genes, including low-abundance species, compared with control samples of a similar sequencing depth. MEM enabled the culture-independent assembly of whole microbial genomes at relative abundances as low as 1% directly from human intestinal biopsies. The assembly of MAGs enables investigation of subpopulation variation across individuals and within an individual's GI tract.

We acknowledge the following limitations of the MEM approach. We have analyzed biopsies with as few as 10<sup>2</sup> 16S rRNA gene copies per  $\mu$ l in the 100  $\mu$ l of elution from the extraction column (corresponding to roughly 10<sup>4</sup> 16S copies per mg of tissue), however, deep analysis of

samples below this bacterial load will require greater levels of host depletion and/or greater sequencing depth (Extended Data Figs. 1 and 2). We advise users to refer to Extended Data Fig. 1 to predict the percentage of nonhost reads from bacterial load to guide sequencing depth decision. We have successfully applied MEM to healthy intestinal biopsies but additional validation should be performed on samples with characteristics that interfere with analysis, for example samples with active inflammation or bleeding. We have successfully applied MEM to fresh samples, but additional validation are needed for preserved tissue samples. We have only characterized the impact of MEM on bacteria and archaea. Future studies will illuminate whether MEM affects the mycobiome and virome.

Here, MEM was validated on mouse feces and intestinal scrapings, rat colon sections, human saliva and human intestinal biopsies. To extend the use of MEM beyond mammals, future studies will be needed to optimize and validate MEM on samples from plants<sup>43,44</sup>, insects<sup>45,46</sup> and other nonmammalian hosts. Sample processing with MEM will also enable higher throughput and less expensive microbiome investigations even in samples with moderate host loads (for example, saliva in which enrichment of microbial reads from 10 to 95% cuts sequencing costs by an order of magnitude). MEM would also benefit researchers investigating evolution and dynamics of microbes and microbial genes across time and across interconnected ecological niches, such as within the human GI tract. In clinical studies, we anticipate that MEM will provide researchers the capability to investigate tumor microbiomes<sup>2,3,47,48</sup>, mucosal intestinal microbiomes<sup>4,9,12,49–51</sup>, tissue translocation of gut microbes<sup>52,53</sup> and the roles of tissue-associated microbes in complex immune disorders<sup>9,54,55</sup>, immune modulation<sup>56</sup> and cancer development<sup>57,58</sup>.

### Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41592-023-02025-4>.

### References

1. Tuganbaev, T. et al. Diet diurnally regulates small intestinal microbiome-epithelial-immune homeostasis and enteritis. *Cell* **182**, 1441–1459 (2020).
2. Dejea, C. M. et al. Patients with familial adenomatous polyposis harbor colonic biofilms containing tumorigenic bacteria. *Science* **359**, 592–597 (2018).
3. Bullman, S. et al. Analysis of *Fusobacterium* persistence and antibiotic response in colorectal cancer. *Science* **358**, 1443–1448 (2017).
4. Morgan, X. C. et al. Dysfunction of the intestinal microbiome in inflammatory bowel disease and treatment. *Genome Biol.* **13**, R79 (2012).
5. Caruso, R., Lo, B. C. & Nunez, G. Host-microbiota interactions in inflammatory bowel disease. *Nat. Rev. Immunol.* **20**, 411–426 (2020).
6. Pascal, V. et al. A microbial signature for Crohn's disease. *Gut* **66**, 813–822 (2017).
7. Gevers, D. et al. The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host Microbe* **15**, 382–392 (2014).
8. Cheng, J. et al. Duodenal microbiota composition and mucosal homeostasis in pediatric celiac disease. *BMC Gastroenterol.* **13**, 113 (2013).
9. Earley, Z. M. et al. GATA4 controls regionalization of tissue immunity and commensal-driven immunopathology. *Immunity* **56**, 43–57 (2023).
10. Ringel, Y. et al. High throughput sequencing reveals distinct microbial populations within the mucosal and luminal niches in healthy individuals. *Gut Microbes* **6**, 173–181 (2015).

11. Parthasarathy, G. et al. Relationship between microbiota of the colonic mucosa vs feces and symptoms, colonic transit, and methane production in female patients with chronic constipation. *Gastroenterology* **150**, 367–379 (2016).
12. Vaga, S. et al. Compositional and functional differences of the mucosal microbiota along the intestine of healthy individuals. *Sci. Rep.* **10**, 14977 (2020).
13. Shen, T. D. et al. The mucosally-adherent rectal microbiota contains features unique to alcohol-related cirrhosis. *Gut Microbes* **13**, 1987781 (2021).
14. Klindworth, A. et al. Evaluation of general 16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res.* **41**, e1 (2013).
15. Chen, L. X., Anantharaman, K., Shaiber, A., Eren, A. M. & Banfield, J. F. Accurate and complete genomes from metagenomes. *Genome Res.* **30**, 315–333 (2020).
16. Vineis, J. H. et al. Patient-specific *Bacteroides* genome variants in pouchitis. *mBio.* **7**, 10–1128 (2016).
17. Groussin, M., Mazel, F. & Alm, E. J. Co-evolution and co-speciation of host-gut bacteria systems. *Cell Host Microbe* **28**, 12–22 (2020).
18. Wang, G. H., Dittmer, J., Douglas, B., Huang, L. & Brucker, R. M. Coadaptation between host genome and microbiome under long-term xenobiotic-induced selection. *Sci. Adv.* **7**, eabd4473 (2021).
19. Tyson, G. W. et al. Community structure and metabolism through reconstruction of microbial genomes from the environment. *Nature* **428**, 37–43 (2004).
20. Pereira-Marques, J. et al. Impact of host DNA and sequencing depth on the taxonomic resolution of whole metagenome sequencing for microbiome analysis. *Front. Microbiol.* **10**, 1277 (2019).
21. Marotz, C. A. et al. Improving saliva shotgun metagenomics by chemical host DNA depletion. *Microbiome* **6**, 42 (2018).
22. Bruggeling, C. E. et al. Optimized bacterial DNA isolation method for microbiome analysis of human tissues. *Microbiology Open* **10**, e1191 (2021).
23. Ganda, E. et al. DNA extraction and host depletion methods significantly impact and potentially bias bacterial detection in a biological fluid. *mSystems* **6**, e0061921 (2021).
24. Avanzi, C. et al. Red squirrels in the British Isles are infected with leprosy bacilli. *Science* **354**, 744–747 (2016).
25. Charalampous, T. et al. Nanopore metagenomics enables rapid clinical diagnosis of bacterial lower respiratory infection. *Nat. Biotechnol.* **37**, 783–792 (2019).
26. Cheng, W. Y. et al. High sensitivity of shotgun metagenomic sequencing in colon tissue biopsy by host DNA depletion. *Genomics Proteomics Bioinformatics* <https://doi.org/10.1016/j.gpb.2022.09.003> (2022).
27. Oechslin, C. P. et al. Limited correlation of shotgun metagenomics following host depletion and routine diagnostics for viruses and bacteria in low concentrated surrogate and clinical samples. *Front. Cell Infect. Microbiol.* **8**, 375 (2018).
28. Hasan, M. R. et al. Depletion of human DNA in spiked clinical specimens for improvement of sensitivity of pathogen detection by next-generation sequencing. *J. Clin. Microbiol.* **54**, 919–927 (2016).
29. Heravi, F. S., Zakrzewski, M., Vickery, K. & Hu, H. Host DNA depletion efficiency of microbiome DNA enrichment methods in infected tissue samples. *J. Microbiol. Meth.* **170**, 105856 (2020).
30. Shaffer, J. P. et al. A comparison of six DNA extraction protocols for 16S, ITS and shotgun metagenomic sequencing of microbial communities. *BioTechniques* **73**, 34–46 (2022).
31. Hallmaier-Wacker, L. K., Lueert, S., Roos, C. & Knauf, S. The impact of storage buffer, DNA extraction method, and polymerase on microbial analysis. *Sci. Rep.* **8**, 6292 (2018).
32. Bellali, S., Lagier, J. C., Raoult, D. & Bou Khalil, J. Among live and dead bacteria, the optimization of sample collection and processing remains essential in recovering gut microbiota components. *Front. Microbiol.* **10**, 1606 (2019).
33. Barlow, J. T., Bogatyrev, S. R. & Ismagilov, R. F. A quantitative sequencing framework for absolute abundance measurements of mucosal and luminal microbial communities. *Nat. Commun.* **11**, 2590 (2020).
34. Weyrich, L. S. et al. Laboratory contamination over time during low-biomass sample analysis. *Mol. Ecol. Resour.* **19**, 982–996 (2019).
35. Salter, S. J. et al. Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol.* **12**, 87 (2014).
36. Velásquez-Mejía, E. P., de la Cuesta-Zuluaga, J. & Escobar, J. S. Impact of DNA extraction, sample dilution, and reagent contamination on 16S rRNA gene sequencing of human feces. *Appl. Microbiol. Biotechnol.* **102**, 403–411 (2018).
37. Liu, Y., Elworth, R. A. L., Jochum, M. D., Aagaard, K. M. & Treangen, T. J. De novo identification of microbial contaminants in low microbial biomass microbiomes with Squeegee. *Nat. Commun.* **13**, 6799 (2022).
38. Mehta, R. S. et al. Stability of the human faecal microbiome in a cohort of adult men. *Nat. Microbiol.* **3**, 347–355 (2018).
39. Human Microbiome Project, C. Structure, function and diversity of the healthy human microbiome. *Nature* **486**, 207–214 (2012).
40. Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database. *Bioinformatics* **36**, 1925–1927 (2019).
41. Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk v2: memory friendly classification with the Genome Taxonomy Database. *Bioinformatics* **38**, 5315–5316 (2022).
42. Delmont, T. O. & Eren, A. M. Linking pangenomes and metagenomes: the *Prochlorococcus* metapangenome. *PeerJ* **6**, e4320 (2018).
43. Castrillo, G. et al. Root microbiota drive direct integration of phosphate stress and immunity. *Nature* **543**, 513–518 (2017).
44. Fitzpatrick, C. R. et al. Assembly and ecological function of the root microbiome across angiosperm plant species. *Proc. Natl Acad. Sci. USA* **115**, E1157–E1165 (2018).
45. Shin, S. C. et al. *Drosophila* microbiome modulates host developmental and metabolic homeostasis via insulin signaling. *Science* **334**, 670–674 (2011).
46. Motta, E. V. S., Raymann, K. & Moran, N. A. Glyphosate perturbs the gut microbiota of honey bees. *Proc. Natl Acad. Sci. USA* **115**, 10305–10310 (2018).
47. Poore, G. D. et al. Microbiome analyses of blood and tissues suggest cancer diagnostic approach. *Nature* **579**, 567–574 (2020).
48. Riquelme, E. et al. Tumor microbiome diversity and composition influence pancreatic cancer outcomes. *Cell* **178**, 795–806 (2019).
49. Olaisen, M. et al. Bacterial mucosa-associated microbiome in inflamed and proximal noninflamed ileum of patients with Crohn’s disease. *Inflamm. Bowel Dis.* **27**, 12–24 (2021).
50. Liou, M. J. et al. Host cells subdivide nutrient niches into discrete biogeographical microhabitats for gut microbes. *Cell Host Microbe* **30**, 836–847 (2022).
51. Libertucci, J. et al. Inflammation-related differences in mucosa-associated microbiota and intestinal barrier function in colonic Crohn’s disease. *Am. J. Physiol. Gastrointest. Liver Physiol.* **315**, G420–G431 (2018).
52. Brenchley, J. M. & Douek, D. C. Microbial translocation across the GI tract. *Annu. Rev. Immunol.* **30**, 149–173 (2012).

53. Singer, J. R. et al. Preventing dysbiosis of the neonatal mouse intestinal microbiome protects against late-onset sepsis. *Nat. Med* **25**, 1772–1782 (2019).
54. Maynard, C. L., Elson, C. O., Hatton, R. D. & Weaver, C. T. Reciprocal interactions of the intestinal microbiota and immune system. *Nature* **489**, 231–241 (2012).
55. Girdhar, K. et al. A gut microbial peptide and molecular mimicry in the pathogenesis of type 1 diabetes. *Proc. Natl Acad. Sci. USA* **119**, e2120028119 (2022).
56. Gopalakrishnan, V. et al. Gut microbiome modulates response to anti-PD-1 immunotherapy in melanoma patients. *Science* **359**, 97–103 (2018).
57. Ferreira, R. M. et al. Gastric microbial community profiling reveals a dysbiotic cancer-associated microbiota. *Gut* **67**, 226–236 (2018).
58. Wilson, M. R. et al. The human gut bacterial genotoxin colibactin alkylates DNA. *Science* **363**, eaar7785 (2019).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Springer Nature or its licensor (e.g. a society or other partner) holds exclusive rights to this article under a publishing agreement with the author(s) or other rightsholder(s); author self-archiving of the accepted manuscript version of this article is solely governed by the terms of such publishing agreement and applicable law.

© The Author(s), under exclusive licence to Springer Nature America, Inc. 2023

## Methods

### Sample collection

**Mice (stool samples).** All animal husbandry and experiments were approved by the Caltech Institutional Animal Care and Use Committee (IACUC protocol no. 21-1769). Male and female wild-type, nontransgenic surplus mice were used for stool collection. These animals were being fed a standard chow (LabDiet catalog no. 3005740-220) before stool collection. The stool was freshly collected by gently handling the mice. A total of three stool pellets from three different mice were collected at a time and were transferred to clean microfuge tubes with sterile tweezers. Samples were stored on ice for up to 30 min before being processed in the laboratory. A total of 1 ml of saline was added to each stool pellet and the samples were homogenized by pipetting. Homogenized stool samples were diluted threefold in saline and 100  $\mu$ l from each diluted stool sample was processed with various host-depletion methodologies ('MEM' and 'Methodological comparisons with published host-depletion protocols').

**Rat (small intestine and colonic samples).** Tissue collection was performed postmortem through an institutional tissue sharing program that does not require Institutional Animal Care and Use Committee approval. One wild-type Syngap surplus rat was euthanized with CO<sub>2</sub> and the small intestine was removed with sterilized tweezers. The rat was being fed a standard chow (LabDiet catalog no. 3005740-220) but was fasted for 6 h before sample collection.

A portion of the small intestine that appeared clear of content was cut and placed on a petri dish on ice. Any remnant luminal contents were removed by squeezing the intestine with tweezers. The intestine was then cut and opened longitudinally with the mucosa facing upwards. A sterile glass slide was used to scrape the small intestine mucosa and placed into a clean microfuge tube on ice. Samples were stored on ice for up to 30 min before being processed in the laboratory. Mucosal scrapings were mixed and separated into 13 clean microfuge tubes, each tube containing a 2 mg of tissues ('MEM' and 'Methodological comparisons with published host-depletion protocols').

The large intestine was placed on a separate petri dish on ice and any luminal contents were removed by squeezing the intestine with tweezers. The entire large intestine was then cut into 14 evenly sized pieces with a sterile scalpel. Sterile tweezers were used to transfer each intestinal piece into a clean microfuge tube on ice. Samples were stored on ice for up to 30 min before being processed in the laboratory ('MEM' and 'Methodological comparisons with published host-depletion protocols').

**Human (saliva samples).** Human saliva samples were acquired from two healthy adult volunteers and analyzed under California Institute of Technology Institutional Review Board (IRB) protocol no. 21-1092. All participants provided (digital) written informed consent before donation. No personal identifying information was collected at the time of consent and participant specimens were coded. Volunteers were asked not to eat, drink, chew gum, brush their teeth or smoke 30 min before collection. No volunteers had taken systemic antibiotics for at least 2 weeks before donation. Volunteers were instructed to pool saliva in their mouths and spit 2 ml of saliva, ignoring bubbles when estimating volume, into a 15 ml conical tube through a plastic funnel. Before undergoing MEM, saliva samples underwent a DTT pretreatment in some experiments. Saliva was mixed at a 1:1 ratio with fresh DTT (10 mM DTT in 1 $\times$  PBS, Sigma Aldrich catalog no. 43815), vortexed briefly and incubated for 1 min at room temperature before undergoing host-depletion processing ('MEM' and 'Methodological comparisons with published host-depletion protocols').

**Human (tissue samples).** All activities related to enrollment of participants, collection of samples and sample analysis were approved by the University of Chicago IRB and performed under IRB protocol nos.

15573A and 13-1080. Deidentified samples were received at Caltech and analyzed under Caltech IRB protocol no. 21-1083. Adults scheduled for routine colon cancer screenings via colonoscopy at the University of Chicago Medicine were screened for diagnosis and eligibility criteria for enrollment in the study on a weekly basis. Exclusion criteria included: participants with chronic infectious diseases such as human immunodeficiency virus or hepatitis C (HCV); active, untreated *Clostridium difficile* infection; active infection with severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2); intravenous or illicit drug use such as cocaine, heroin, nonprescription methamphetamines; active use of blood thinners; severe comorbid diseases; participants on active cancer treatment and participants who were pregnant. Approaching prospective participants was at the discretion of their treating physician and was not done in cases that would put participants at any increased risk, regardless of reason. Participants were approached the day of their procedure and informed, written consent was obtained before any samples were acquired.

### Human ascending colon paired MEM and control

To assess the impact of MEM on intestinal microbes, eight ascending colon biopsies, designated as 10 cm distal to the ileocecal valve, were collected from a single field of view (5 cm diameter) for five different participants (Fig. 3a). Biopsies were collected in a total of 2–3 passages with 3–4 biopsies per passage using a pair of 2.8-mm biopsy forceps. Biopsies from the same passage were stored together on ice in a dry microfuge tube for an average of 28 min (ranging from 15 to 36 min). After samples were transferred to the laboratory, biopsies from the same passage were split into control and MEM groups for a total of four biopsies per condition with evenly sized biopsies present in each group. Biopsy size ranged from 0.2 to 4.8 mg with an average weight of 2.49 mg. Nonhost-depleted biopsies were processed individually by adding 150  $\mu$ l of PrimeStore MTM inactivation buffer (Longhorn) to each biopsy and vortexing briefly before storing at  $-80^{\circ}\text{C}$  until DNA extraction. Depleted samples were processed individually at University of Chicago ('MEM') before shipment on dry ice to Caltech for DNA extraction.

### Longitudinal sampling of the human intestinal tract

For longitudinal sampling, a total of five participants were sampled 12 times from four different locations during a routine colonoscopy (Fig. 4c). The four locations sampled were the terminal ileum, ascending colon (designated as 10 cm distal to the ileocecal valve), descending colon and rectum. From a single field of view (5 cm diameter) from each location, three biopsies were collected in one passage with 2.8 mm biopsy forceps and stored dry on ice in a microfuge. For participant CT14, only one rectal sample was obtained. On average, biopsies were 2.5 mg with a minimum size of 0.1 mg and a maximum of 5.9 mg. All biopsies were then processed individually in the laboratory at University of Chicago ('MEM') before shipment on dry ice to Caltech for DNA extraction. Time between specimen collection and processing ranged from 10 to 52 min. Samples were processed individually in the laboratory at University of Chicago ('MEM') before shipment on dry ice to Caltech for DNA extraction. Additionally, three microfuge tubes of 400  $\mu$ l of saline were opened and closed in the laboratory and processed with MEM to serve as clinical processing blanks.

### Depletion protocols

**MEM.** Samples for MEM treatment were placed into 2-ml 1.4-mm ceramic bead-beating tubes (Lysing Matrix D from MP Biomedical, catalog no. 116913050-CF) with a maximum volume of 400  $\mu$ l. For solid sample types (stool and intestinal tissue), up to 400  $\mu$ l of saline (0.9% NaCl, autoclaved) was added into the bead-beating tube. Samples were homogenized using FastPrep-24 (MP Biomedical catalog no. 116004500) for 30 s at 4.5 m s<sup>-1</sup> and then immediately placed on ice. A total of 150  $\mu$ l of homogenized tissue was removed and placed into a

clean microfuge tube containing 10  $\mu$ l of buffer (100 mM Tris + 40 mM MgCl<sub>2</sub>, pH 8.0 and 0.22  $\mu$ m sterile filtered), 33  $\mu$ l of saline (0.9% NaCl, autoclaved), 2  $\mu$ l of Benzonase Nuclease HC (EMD Millipore catalog no. 71205) and 5  $\mu$ l of Proteinase K (NEB catalog no. P8107S). Samples were mixed lightly by manually pipetting up and down 5–10 times and spun briefly to pool (1,000g for 5 s). Tubes were placed on a dry block incubator for 15 min at 37 °C with shaking at 600 rpm. Samples were then pelleted at 10,000g for 2 min and the supernatant was removed and discarded. Pellets were resuspended in 150  $\mu$ l of PrimeStore MTM (Longhorn), a transport medium, to inactivate residual enzymatic activity and stored at –80 °C until nucleic-acid extraction. The initial MEM protocol used DNase I treatment in place of Benzonase. However, we noted continuous microbial lysis during DNase I heat inactivation. Benzonase was used to remove high heat steps as it is fully inactivated by PrimeStore MTM.

**Methodological comparisons with published host-depletion protocols.** For all mouse, rat and human saliva samples, the following published protocols were conducted to compare with MEM.

**MolYsis.** Host removal was performed with MolYsis Basic5 (Molzym catalog no. D-301-050) following the manufacturer's protocol. A proteinase K pretreatment (10  $\mu$ l of NEB Proteinase K (catalog no. P8107S)) was performed on solid-tissue samples (stool and intestinal samples) based on Molyzm's recommendations. The entire protocol was performed, including the additional BugLysis step before nucleic-acid extraction ('DNA extraction').

**QIAamp microbiome.** Host removal was performed with QIAamp DNA Microbiome Kit (Qiagen catalog no. 51704) following the manufacturer's protocol. Buffer AHL was aliquoted on kit arrival and was not freeze–thawed more than once. To remove confounding factors from different DNA-extraction kits, the QIAamp DNA Microbiome Kit protocol was followed until the proteinase K incubation and the sample was then processed for nucleic-acid extraction ('DNA extraction').

**lyPMA.** A previously published protocol known as lyPMA, was tested according to the paper's specifications<sup>21</sup>. Liquid samples (diluted stool and saliva) were pelleted at 10,000g for 8 min. Supernatant was removed and the pellet was resuspended in 200  $\mu$ l of nuclease-free water with a light vortex. Samples were left at room temperature for 5 min. After samples were covered with foil, 10  $\mu$ M of PMA (Propidium monoazide) was added and mixed by lightly vortexing each tube for a few seconds. Samples were incubated for 5 min in the dark at room temperature before being placed on ice less than 20 cm from a fluorescent bulb. Samples were incubated under light for 25 min with a quick centrifugation and rotation every 5 min. All samples were then processed for nucleic-acid extraction ('DNA extraction'). The lyPMA method was not tested on rat colonic sectionals due to the limited efficacy of osmotic lysis on solid tissues seen from mouse mucosa samples.

**DNA extraction.** Nucleic acids were isolated following Qiagen's AllPrep PowerViral DNA/RNA Kit (catalog no. 28000-50). Samples were homogenized in 0.1 mm glass beads for 1 min at 6 m s<sup>-1</sup> using FastPrep-24 (MP Biomedical catalog no. 116004500) to ensure complete microbial lysis<sup>59</sup>. A maximum of 24 clinical samples were processed at a time and at least three processing blanks were run on each extraction kit. Samples were eluted into 100  $\mu$ l of nuclease-free water. It should be noted that standard microbial bead beating with 0.1 mm beads was not sufficient to completely lyse intact (control) biopsies in this study. Control biopsies were homogenized with Lysis Matrix E (MP Biomedical catalog no. 116914050-CF) for 1 min at 6 m s<sup>-1</sup> three times with 5 min of incubation on ice between each bead beating.

**Quantification of host DNA.** Host load present in extracted DNA was characterized by droplet digital PCR (ddPCR) of a single-copy gene. For human saliva and tissue samples, the gene *EIF5B* was amplified based on primers found from literature<sup>60</sup> (Forward: 5'-GCCAACTTCAGCCTTCTCTTC-3' and Reverse: 5'-CTCTGGCAACATTTCACACTACA-3'). For samples originating from rodents, the gene *Cyp8b1* was amplified based on primers found from literature<sup>61</sup> (Forward: 5'-GGCTGGCTTCTGAGCTTATT-3' and Reverse: 5'-ACTTCTGAACAGCTCATCGG-3'). Samples were amplified on the C-1000 thermocycler (Bio-Rad, catalog no. 1851196) and quantified using the QX200 ddPCR system (Bio-Rad, catalog no. 1864001). The concentrations of the components in the ddPCR mix used in this study were as follows: 1 $\times$  QX200 ddPCR EvaGreen SuperMix (Bio-Rad, catalog no. 1864035), 500 nM forward primer and 500 nM reverse primer for a total reaction volume of 25  $\mu$ l. Thermocycling was performed as follows: 95 °C for 5 min, 40 cycles of 95 °C for 30 s, 60 °C for 30 s and 68 °C for 30 s, followed by a dye-stabilization step at 4 °C for 5 min and 90 °C for 5 min. All ramp rates were 2 °C s<sup>-1</sup>. LOB refers to limit of blank defined as LoB = meanblank + 1.645[SDBblank] based on three processing blanks.

**Quant-seq.** Microbial characterization and quantification were performed using the quantitative sequencing (Quant-seq) pipeline we have described previously<sup>33</sup>. Due to the low bacterial loads present in intestinal biopsies, Quant-seq was also performed on three MEM processing blanks. If a taxon was detected at a higher absolute abundance in any of the processing blanks compared to the intestinal biopsies, the taxon was removed from downstream analysis. Only genera and species above 0.1% abundance in at least one biopsy were considered for the analysis in Figs. 4d–f.

**Shotgun sequencing.** Extracted DNA was prepared for sequencing using Illumina DNA Prep (catalog no. 20018704). A maximum input of 500 ng of DNA was used for library prep. After processing with the MEM protocol, almost all human biopsy samples had less than Illumina's recommended minimal DNA input amount of 1 ng and were below the limit of detection of the Qubit double-stranded DNA (dsDNA) High Sensitivity assay (Thermo catalog no. Q32851). Estimations of input DNA were made using 16S rRNA gene ddPCR ('Quant-seq') and host quantification ('Quantification of host DNA' and equation (1)). For these calculations, we assumed the 16S rRNA gene copy number (four per cell), total DNA per microbial cell (3fg based on average genome size of 3 Mb) and absence of nonhost and/or nonprokaryotic DNA.

$$30\mu\text{l} \times \text{prokaryotic load} \left( \frac{16\text{S rRNA gene copies}}{\mu\text{l}} \right) \times \frac{1\text{ cell}}{4 \times 16\text{S rRNA gene copies}} \times \frac{3\text{fg DNA}}{1\text{ cell}} = \text{prokaryotic DNA (fg)} \quad (1)$$

$$\text{prokaryotic DNA} + \text{host DNA} = \text{total DNA.}$$

For samples with DNA concentrations below Illumina's recommended input, additional PCR cycles were added to the amplification step based on DNA input (Supplementary Table 11).

Finished libraries were quantified through Qubit's dsDNA High Sensitivity assay and a High Sensitivity D1000 TapeStation Chip (Agilent catalog nos. 5067-5585, 5067-5584). If additional peaks were seen at 45 or 120 bp, indicating the presence of primer dimers or adapter dimers, we performed an additional clean-up step with AMPureXP beads (Beckman Coulter, catalog no. A63880) at a ratio of 0.8:1 of beads to library volume. For quantification, finalized libraries were amplified on the CFX-96 quantitative PCR (qPCR) (Bio-Rad catalog no. 1855196) with primers targeting the Illumina adapter sequence (Forward: 5'-AAT GAT ACG GCG ACC ACC GA-3' and Reverse: 5'-CAA GCA GAA GAC GGC ATA CGA-3'). Libraries were diluted 1:40,000 in nuclease-free water before amplification to fall within the range of KAPA standards concentrations (Roche catalog

no. 07960387001) for quantification. The concentrations of the components in the qPCR mix used were as follows: 1× SsoFast EvaGreen Supermix (Bio-Rad catalog no. 1725201), 125 nM forward primer and 125 nM reverse primer for a total reaction volume of 10 µl. Thermocycling was performed as follows: 95 °C for 5 min, 40 cycles of 95 °C for 30 s and 60 °C for 45 s, followed by a melt-curve step at 95 °C for 15 s, 50 °C for 15 s, 70 °C for 1 s and 95 °C for 5 s. Pooled samples were quantified through Qubit's dsDNA High Sensitivity assay and a High Sensitivity D1000 TapeStation Chip before submitting the samples for sequencing. Sequencing was performed by Fulgent Genetics using the Illumina NovaSeq6000 platform. Sequencing batch 1 was performed on the NovaSeq6000 SP flow cell and 2 × 100 bp reagent kit for paired-end sequencing with an average sequencing depth of 23 million reads. Sequencing batch 2 was used for MAG assembly and was performed on one NovaSeq6000 S4 lane and 2 × 150 bp reagent kit for paired-end sequencing with an average sequencing depth of 223 million reads. Samples were demultiplexed on the NovaSeq6000 and raw fastq files for read 1 and read 2 were provided along with fastqc files for each sample.

The number of nonhost reads obtained from each sample can be accurately predicted based on a single qPCR measurement of bacterial load (16S rRNA gene copies) (Extended Data Fig. 1) and can be used to inform necessary sequencing depth.

**Marker gene analyses.** Sequencing data were processed using the KneadData v.0.10.0 (ref. 62). Through KneadData, quality control and host removal were performed with Trimmomatic v.0.39 (refs. 63). Human derived sample types were aligned to KneadData's default human reference genome (a combination of hg38 human genome reference (GenBank assembly accession no. [GCA\\_000001405.29](#)) and small contaminant sequences) and aligned reads were removed. Samples acquired from mice were processed using the reference genome GRCm39 constructed from C57BL/6J mouse strains (GenBank assembly accession no. [GCA\\_000001635.9](#)). After bioinformatic host removal, the percentages of host reads were calculated by dividing reads remaining after host filtering by the total reads that passed quality control. To assign species, nonhost reads from read 1 and read 2 were then concatenated and processed using the MetaPhlan v.3.0 workflow outlined in bioBakery (<https://github.com/biobakery/biobakery>) under default settings (Database mpa\_v30\_CHOCOPHlan\_201901)<sup>62</sup>. For stool, nearly 90% of the nonhost reads did not align to known bacteria in the HUMAnN databases, likely due to the bias toward human microbiome datasets.

**HUMAnN pathway and gene alignment.** Nonhost read 1 and read 2 outputted from KneadData were concatenated and processed using the HUMAnN v.3.0 workflow outlined in bioBakery (<https://github.com/biobakery/biobakery>) under default settings<sup>62</sup>. Taxonomic profiles obtained from MetaPhlan ('Marker gene analyses') were merged within participants and used as taxonomic inputs using the '-taxonomic-profile' flag in HUMAnN. Reported pathway abundances and gene abundances were normalized to relative abundances and concatenated.

**MAG assembly.** Sequencing data were processed using the metagenomic workflow<sup>64,65</sup> outlined in anvi'o<sup>66,67</sup> v.7.1 (<https://anvio.org>). Quality control filtering of short reads was performed using the Illumina-utils library<sup>68</sup> v.2.12. Host reads were removed by alignment to the hg38 human genome reference (GenBank assembly accession no. [GCA\\_000001405.29](#)). Assembly was performed on each sample individually using MEGAHIT<sup>69</sup> v.1.2.9 unless coassembly was explicitly stated as in Fig. 4, with default setting except setting a minimum contig length of 1,000 bp. Short reads generated from each sample were then aligned to contigs generated from all assemblies using Bowtie2 (refs. 70) v.2.3.5. Contigs were processed using anvi'o to generate a

contig databases with the command 'anvi-gen-contigs-database' with default settings and with Prodigal<sup>71</sup> v.2.6.3 to identify open reading frames. Single-copy core genes were detected with 'anvi-run-hmm' to (bacteria  $n = 71$  and archaea  $n = 76$ , modified from Lee in ref. 72, rRNAs ( $n = 12$ , modified from <https://github.com/tseemann/bar-rnap>) using HMMer<sup>73,74</sup> v.3.3.2. Genes were annotated using both 'anvi-run-ncbi-cogs' for the National Center for Biotechnology Information (NCBI)'s COGs database<sup>75</sup> and 'anvi-run-kegg-kofams' from the Kofam HMM (hidden Markov model) database of Kyoto Encyclopedia of Genes and Genomes orthologs (knock-outs)<sup>76</sup>. BAM files were profiled with 'anvi-profile' and merged with 'anvi-merge' for samples originating from the same participant. Automatic binning was performed by CONCOCT<sup>77</sup> v.1.1.0 by specifying a maximum number of bins based on the estimated number of bacterial genomes computed from each sample's contigs. The maximum number of bins was set to one out of three of the number of expected genomes to limit the likelihood of fragmentation. Bins generated with CONCOCT were imported in the anvi'o profile database and were then manually refined and summarized to obtain fasta files of individual MAGs. Once manual binning of all samples from the same participant was complete, MAGs above 50% complete were dereplicated to generate a unique list of genomes using anvi'o and pyani v.0.2.11. Representative genomes were chosen based on quality scores and clustered based on more than 95% ANI. The final list of MAGs was taxonomically assigned with GTDB-Tk (Genome Taxonomy Database Toolkit; v.2.1.0, refs. 40,41) using classify\_wf with default settings.

**Strain analysis across individuals.** A *P. vulgatus* MAG from the terminal ileum of CT12 was selected as a reference genome based on genome length. Open reading frames were identified through Prodigal for the *P. vulgatus* reference genome. Nonhost reads from each participant (CT7, CT12, CT13 and CT14) were mapped onto the *P. vulgatus* reference genome by following anvi'o's metagenomics workflow using reference mode. For each sample and each gene present in the *P. vulgatus* reference genome, gene detection was calculated. Gene detection refers to the percentage of each gene sequence with at least 1× coverage. The average detection across all genes present within the *P. vulgatus* MAG was calculated and samples with a mean detection below 0.25 were removed from the final analysis. Pangenome visualization was performed in anvi'o interactive interface using the gene-mode flag with sorting of samples and genes by detection.

**Analyses of SNVs.** One terminal ileum sample from CT12 was split into three technical replicates before library preparation, and each replicate was sequenced at a depth of 150 million to 250 million reads in sequencing batch 2. SNV analyses across these samples were performed with anvi'o after dereplication ('MAG assembly') using the command 'anvi-gen-variability-profile' with a minimum mean coverage of 50× in all samples. Biological SNVs were classified as being present in all three technical replicates. SNVs present in only one or two technical replicates were classified as sequencing, PCR or input errors. A threshold for minimum deviation from consensus was set based on the deviation required for all SNVs to be present in all technical replicates. This analysis was repeated for each MAG of interest (min mean coverage of 50×,  $n = 6$ ). After a threshold for minimum deviation from consensus was established, longitudinal samples from participant CT12 were analyzed using 'anvi-gen-variability-profile' at the nucleotide, codon and amino acid level with the same minimum mean coverage of 50× and filtering out SNVs occurring in only one sample. The fixation index was computed using 'anvi-gen-fixation-index-matrix' to describe the population structure between samples.

## Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

## Data availability

The datasets generated and analyzed during the current study are available at CaltechDATA, <https://doi.org/10.22002/gx69z-wec80>. Microbial sequencing data are available at NCBI Accession no. [PRJNA991155](https://www.ncbi.nlm.nih.gov/PRJNA991155). Sequencing data from human samples have been host scrubbed using STAT<sup>78</sup> sra-human-scrubber (<https://github.com/ncbi/sra-human-scrubber>) followed by alignment to CHM13 (ref. 79). Source data are provided with this paper.

## Code availability

The code used in data processing and analysis is available at CaltechDATA, <https://doi.org/10.22002/gx69z-wec80>.

## References

59. Mancabelli, L. et al. Guideline for the analysis of the microbial communities of the human upper airways. *J. Oral. Microbiol.* **14**, 2103282 (2022).
60. Kline, M. C., Romsos, E. L. & Duewer, D. L. Evaluating digital PCR for the quantification of human genomic DNA: accessible amplifiable targets. *Anal. Chem.* **88**, 2132–2139 (2016).
61. Zhang, X., Osaka, T. & Tsuneda, S. Bacterial metabolites directly modulate farnesoid X receptor activity. *Nutr. Metab.* **12**, 48 (2015).
62. Beghini, F. et al. Integrating taxonomic, functional, and strain-level profiling of diverse microbial communities with bioBakery 3. *eLife* **10**, e65088 (2021).
63. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
64. Shaiber, A. et al. Functional and genetic markers of niche partitioning among enigmatic members of the human oral microbiome. *Genome Biol.* **21**, 292 (2020).
65. Koster, J. & Rahmann, S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012).
66. Eren, A. M. et al. Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**, e1319 (2015).
67. Eren, A. M. et al. Community-led, integrated, reproducible multi-omics with anvi'o. *Nat. Microbiol.* **6**, 3–6 (2021).
68. Eren, A. M., Vineis, J. H., Morrison, H. G. & Sogin, M. L. A filtering method to generate high quality short reads using illumina paired-end technology. *PLoS ONE* **8**, e66643 (2013).
69. Li, D., Liu, C. M., Luo, R., Sadakane, K. & Lam, T. W. MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics* **31**, 1674–1676 (2015).
70. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
71. Hyatt, D. et al. Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **11**, 119 (2010).
72. Lee, M. D. GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics* **35**, 4162–4164 (2019).
73. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res.* **39**, W29–W37 (2011).
74. Eddy, S. R. Accelerated profile HMM searches. *PLoS Comput Biol.* **7**, e1002195 (2011).
75. Tatusov, R. L. et al. The COG database: an updated version includes eukaryotes. *BMC Bioinformatics* **4**, 41 (2003).
76. Kanehisa, M. & Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes. *Nucleic Acids Res.* **28**, 27–30 (2000).
77. Alneberg, J. et al. Binning metagenomic contigs by coverage and composition. *Nat. Methods* **11**, 1144–1146 (2014).
78. Katz, K. S. et al. STAT: a fast, scalable, MinHash-based k-mer tool to assess Sequence Read Archive next-generation sequence submissions. *Genome Biol.* **22**, 270 (2021).

79. Nurk, S. et al. The complete sequence of a human genome. *Science* **376**, 44–53 (2022).

## Acknowledgements

We acknowledge assistance with animal experiments from Caltech Office of Laboratory Animal Research. We thank M. Ratanapanichkich (California Institute of Technology) for assistance on manual refinement of metagenomic bins and feedback on figure design. We thank A. Carter (California Institute of Technology) for assistance with Quant-seq library preparation, ddPCR measurements and feedback during manuscript preparation. We thank M. Cooper (California Institute of Technology) for identifying appropriate statistical tests, guidance during Quant-seq analysis and feedback on figure design. We thank S. R. Bogatyrev for preliminary investigations, discussions and advice. We thank O. Pradhan (California Institute of Technology) and R. Akana (California Institute of Technology) for advice and feedback during manuscript preparation. We thank B. McDonald (University of Chicago) for providing his expertise and advice on clinical sample collection and processing. We thank A. Wang (University of Chicago) for her assistance in the processing of the human tissue for Figs. 3–6. We thank N. Shelby (California Institute of Technology) for contributions to writing and editing this manuscript. This work was funded in part by a grant from the Kenneth Rainin Foundation (grant no. 2018-1207 to R.F.I.), the Army Research Office Multidisciplinary University Research Initiative (grant no. W911NF-17-1-0402 to R.F.I.), the Jacobs Institute for Molecular Engineering for Medicine, a NIH NIDDK grant (no. RC2 DK133947 to R.F.I. and B.J.), a National Science Foundation Graduate Research Fellowship (grant no. DGE-1745301 to N.J.W.-W.), and a National Institutes of Health Biotechnology Leadership Pre-doctoral Training Program fellowship from Caltech's Donna and Benjamin M. Rosen Bioengineering Center (grant no. T32GM112592, to J.T.B.), a Helmsley Foundation grant (to F.T.), a NIH NIDDK grant (no. RC2 DK122394, to F.T.), a F30 (grant no. 5F30DK121470, to D.G.S.), a R01 (grant no. DK067180, to B.J.) and the Digestive Diseases Research Core Center grant no. P30 DK42086 at the University of Chicago (to B.J.). The funders had no role in the design of the study, the collection, analysis and interpretation of data, nor in writing the manuscript.

## Author contributions

N.J.W.-W. and J.T.B. conceived and optimized MEM. J.T.B. designed sample collection and analyzed 16S sequencing. D.G.S. codesigned and performed human biopsy collection. N.J.W.-W. and F.T. analyzed shotgun sequencing. A.E.R. performed library preparation. R.F.I. contributed to the design and implementation of the study and to obtaining funding. A.M.E. oversaw the bioinformatic analysis, contributed to the design and implementation of the study and to obtaining funding. B.J. supervised the clinical work, contributed to the design and implementation of the study and to obtaining funding. All authors edited the manuscript. A detailed author contribution statement is available in the Supplementary Information.

## Competing interests

The work in this paper is the subject of a patent application filed by Caltech (R.F.I., N.J.W.-W., J.T.B. and A.E.R.). The other authors declare no competing interests.

## Additional information

**Extended data** is available for this paper at <https://doi.org/10.1038/s41592-023-02025-4>.

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s41592-023-02025-4>.

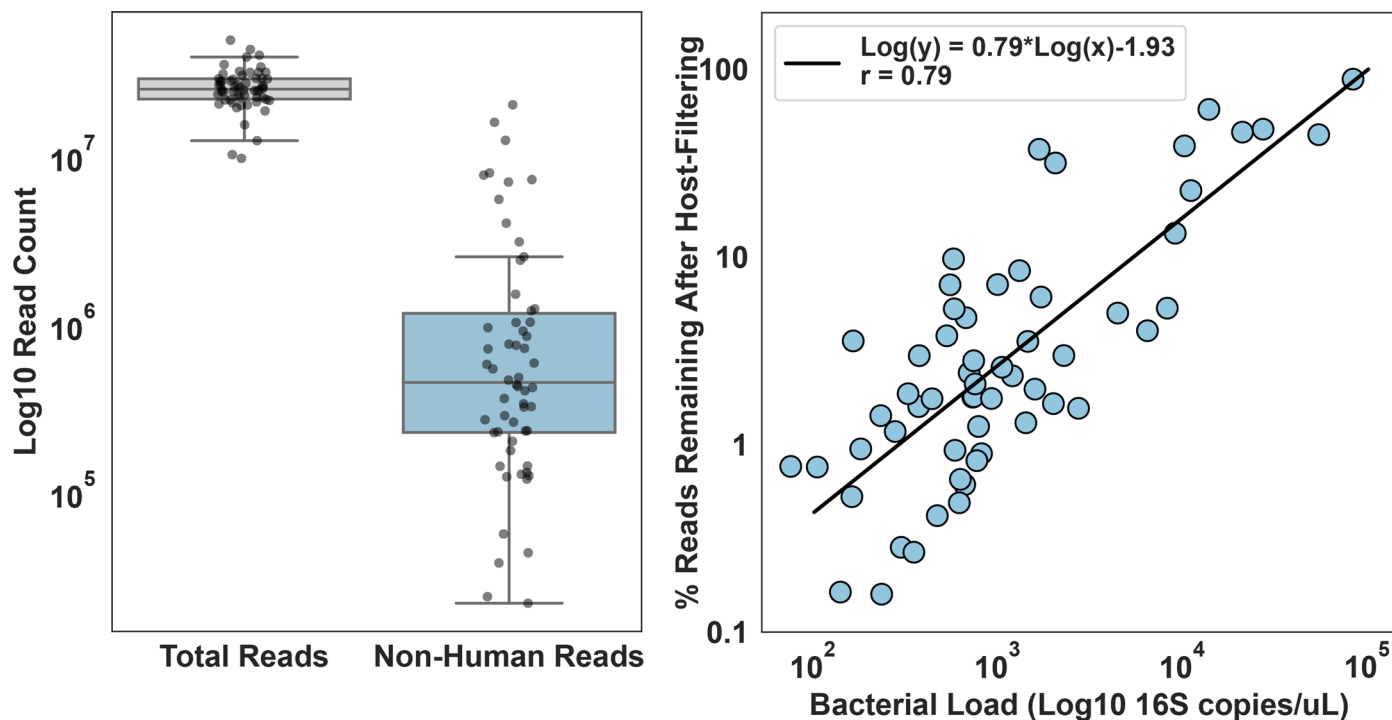
**Correspondence and requests for materials** should be addressed to Rustem F. Ismagilov.

Peer reviewer reports are available. Primary Handling Editor: Lei Tang and Hui Hua, in collaboration with the *Nature Methods* team.

**Peer review information** *Nature Methods* thanks the anonymous reviewers for their contribution to the peer review of this work.

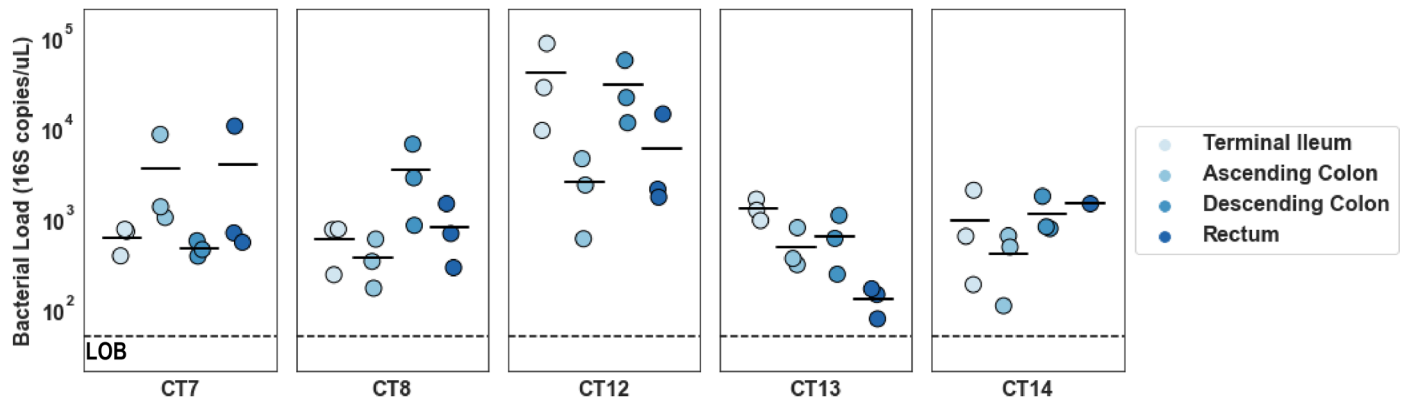
**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).



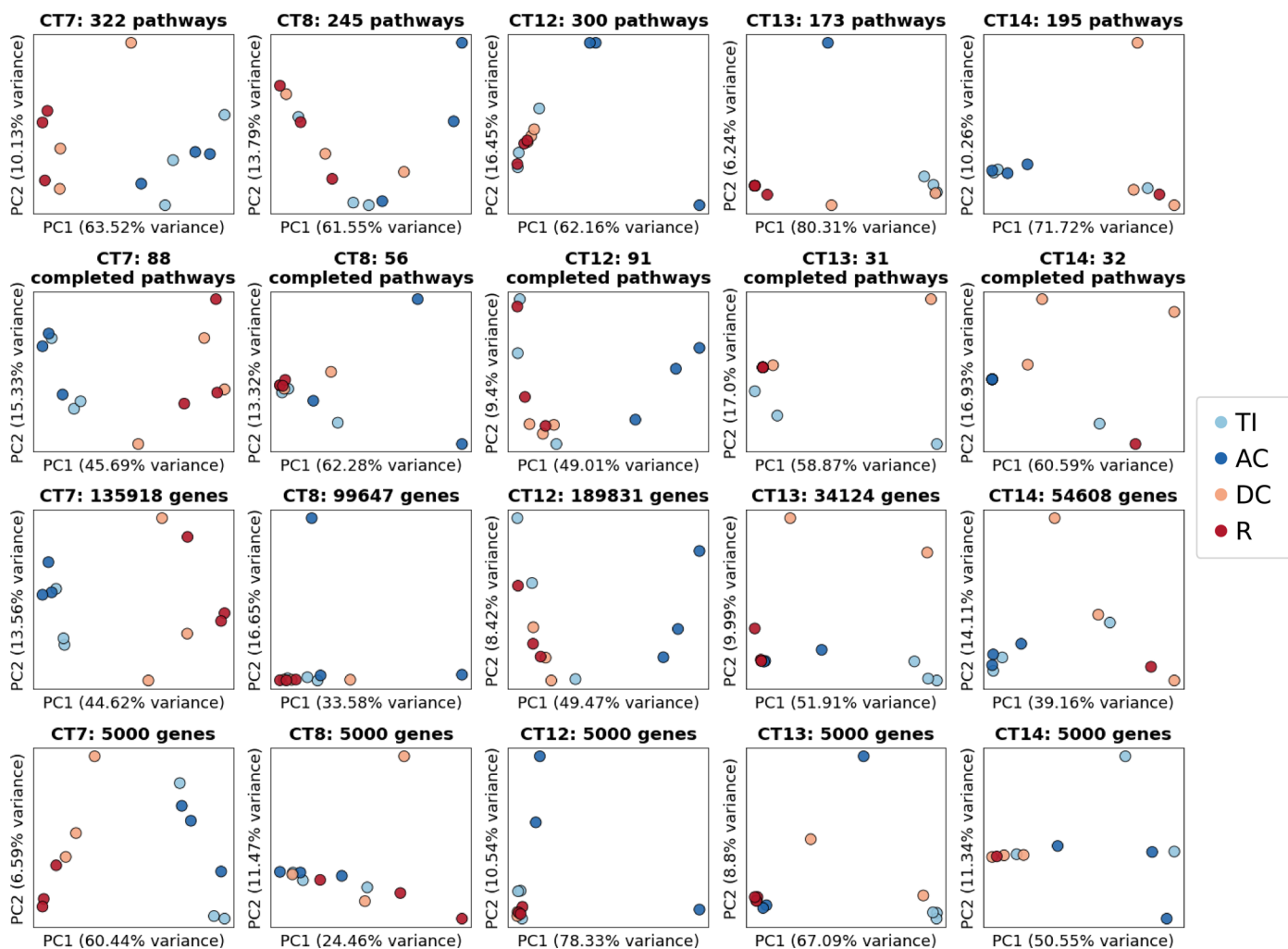


**Extended Data Fig. 1 | Correlation between bacterial load and non-host reads.** Shotgun sequencing was performed on longitudinally sampled intestinal biopsies after processing with host depletion (N = 60 biological replicates). Roughly 25 million reads on average were obtained for each biopsy and all samples fit on a single NovaSeq S1 flowcell. After host-filtering an average of 2 million reads were remaining with a range from 2E4 reads to 2E7 reads. For each box, the middle horizontal line denotes mean values, boxes extend to the 25th

and 75th percentile, and whiskers extend to the 1.5 interquartile range. The variability in non-host reads remaining had a strong correlation (Spearman,  $r = 0.79$ ) with the total microbial load as measured by digital PCR. This strong correlation indicated that our process was achieving a relatively uniform depletion across all samples. Additionally, the strong correlation indicates that the majority of non-human reads in our samples come from bacteria picked up by the 16S primers used for total microbial load quantification.

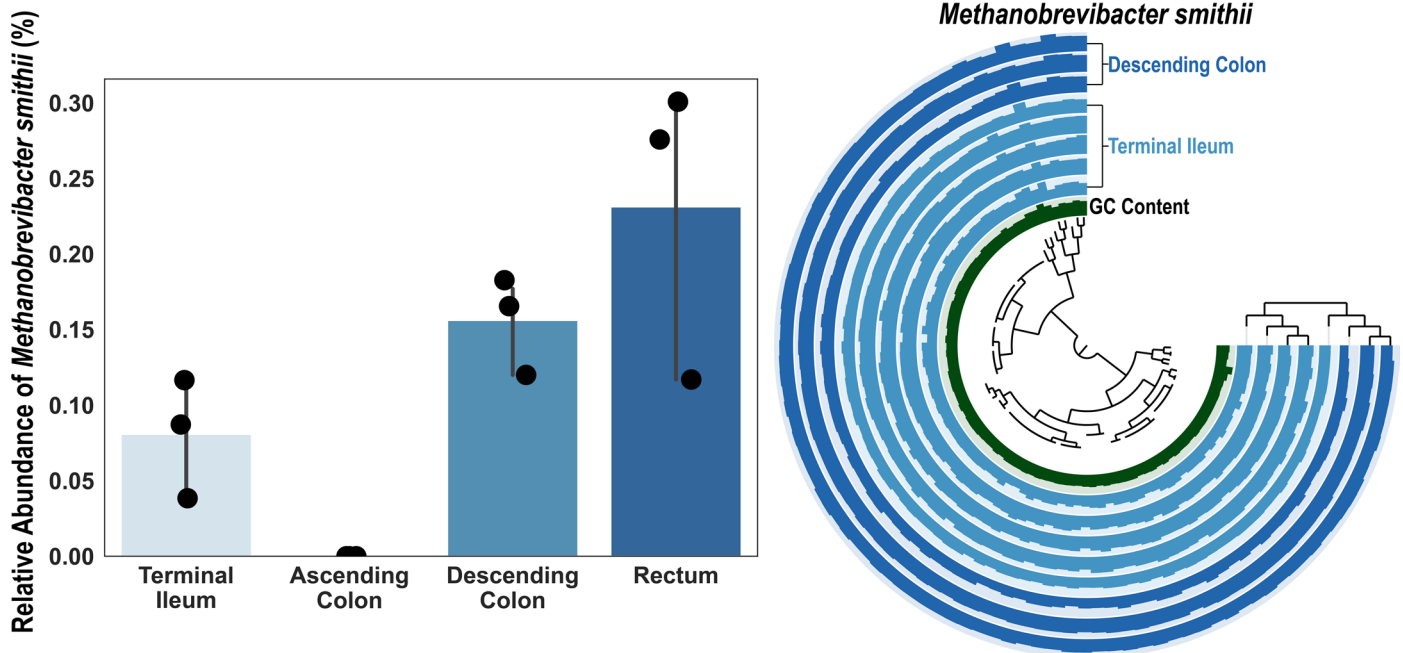


**Extended Data Fig. 2 | Bacterial loads of longitudinal biopsies.** 16S rRNA gene copies were quantified as a proxy for bacterial load for all biopsies. Samples were plotted by participant and then by location. (N = 3 biological replicates for each location for each participant, LOB refers to limit of blank defined as  $LoB = \text{mean}_{\text{blank}} + 1.645[SD_{\text{blank}}]$  based on three processing blanks).



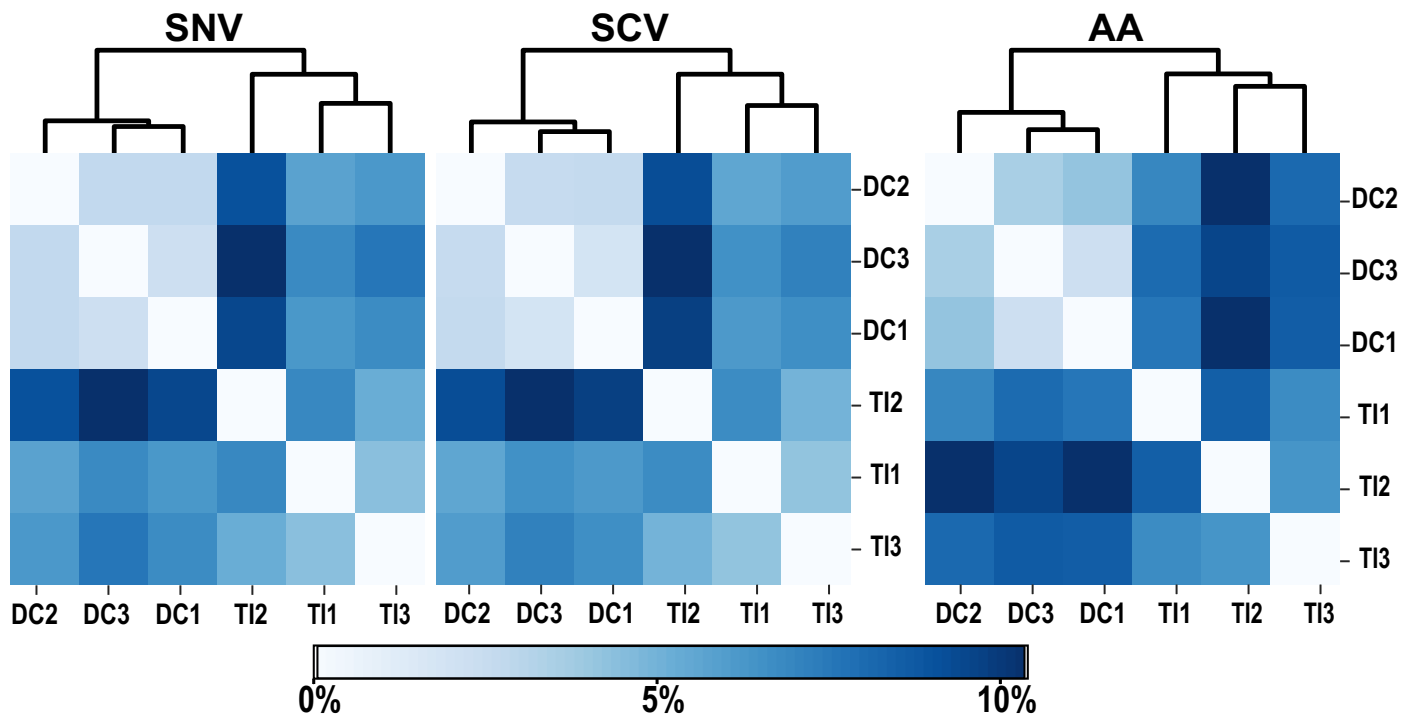
**Extended Data Fig. 3 | Longitudinal variation at the pathways and gene-level.** PCA analysis was performed on all 60 longitudinal samples grouped by participant (CT7, CT8, CT12, CT13, and CT14). Shotgun-sequencing data was annotated for pathways and genes through HUMAnN 3 without the taxonomic-

profile flag. A) PCA on relative abundance of all pathways. B) PCA on relative abundance of completed pathways (defined as above 90% of modules being present). C) PCA on relative abundance of all genes. D) PCA on relative abundance of the top 5,000 most abundant genes in each participant.



**Extended Data Fig. 4 | Archaeon *Methanobrevibacter smithii* found along the lower GI tract.** From shotgun sequencing, we detected participant CT12 had low levels of *Methanobrevibacter smithii* present in the terminal ileum, descending colon, and rectal biopsies (N = 3 biological replicates; error bars are 95% CI

centered on the mean). MAG construction was performed on co-assembly of all biopsies taken from the terminal ileum and descending colon to reconstruct a full *Methanobrevibacter smithii* genome (completeness: 100%, redundancy: 0%).



## Reporting Summary

Nature Research wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Research policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

### Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

- |                                     |                                                                                                                                                                                                                                                                                                |
|-------------------------------------|------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| n/a                                 | Confirmed                                                                                                                                                                                                                                                                                      |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> The exact sample size ( $n$ ) for each experimental group/condition, given as a discrete number and unit of measurement                                                                                                                                    |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly                                                                                                                                    |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> The statistical test(s) used AND whether they are one- or two-sided<br><i>Only common tests should be described solely by name; describe more complex techniques in the Methods section.</i>                                                               |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> A description of all covariates tested                                                                                                                                                                                                                     |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons                                                                                                                                        |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals) |
| <input checked="" type="checkbox"/> | <input checked="" type="checkbox"/> For null hypothesis testing, the test statistic (e.g. $F$ , $t$ , $r$ ) with confidence intervals, effect sizes, degrees of freedom and $P$ value noted<br><i>Give <math>P</math> values as exact values whenever suitable.</i>                            |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings                                                                                                                                                                      |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes                                                                                                                                                |
| <input checked="" type="checkbox"/> | <input type="checkbox"/> Estimates of effect sizes (e.g. Cohen's $d$ , Pearson's $r$ ), indicating how they were calculated                                                                                                                                                                    |

*Our web collection on [statistics for biologists](#) contains articles on many of the points above.*

### Software and code

Policy information about [availability of computer code](#)

Data collection

Data analysis https://github.com/biobakery/biobakery) under default settings (Database: mpa\_v30\_CHOCOPhlan\_201901). Non-host read 1 and read 2 outputted from KneadData were concatenated and processed using the HUMAnN 3.0 workflow outlined in bioBakery (<https://github.com/biobakery/biobakery>) under default settings. Taxonomic profiles obtained from MetaPhlan (see "Marker Gene Analyses") were merged within patients and used as taxonomic inputs using the "--taxonomic-profile" flag in HUMAnN. Reported pathway abundances and gene abundances were normalized to relative abundances and concatenated. Sequencing data was processed using the metagenomic workflow outlined in anvio v7.1 (<https://anvio.org>). QC filtering of short reads was performed using the Illumina-utils library v2.12. Assembly was performed on each sample individually using MEGAHIT v1.2.9 unless co-assembly was explicitly stated as in figure 4, with default setting except setting a minimum contig length of 1000bp. Short reads generated from each sample were then aligned to contigs generated from all assemblies using Bowtie2 v2.3.5. Contigs were processed using anvio to generate a contig databases with the command "anvi-gen-contigs-database" with default settings and with Prodigal v2.6.3 to identify open reading frames. Single-copy core genes were detected with "anvi-run-hmm" to (bacteria  $n = 71$  and archaea  $n = 76$ , modified from Lee, 2019, ribosomal RNAs (rRNAs) ( $n = 12$ , modified from <https://github.com/tseemann/barrnap>) using HMMer v3.3.2. Genes were annotated using "anvi-run-ncbi-cogs" for NCBI's Clusters of Orthologous Groups (COGs) database and "anvi-run-kegg-kofams" from the KOFam HMM database of KEGG orthologs (KOs). BAM files were profiled with "anvi-profile" and merged with "anvi-merge" for samples originating from the same participant. Automatic binning was performed by CONCOCT v1.1.0 by specifying a maximum number of bins based on the estimated number of bacterial genomes computed from each sample's contigs. The maximum number of bins was set to 1/3 the number of expected

genomes to limit the likelihood of fragmentation. Bins generated with CONCOCT were imported in the anvi'o profile database and were then manually refined and summarized to obtain fasta files of individual MAGs. Once manual binning of all samples from the same participant was complete, MAGs above 50% complete were dereplicated to generate a unique list of genomes using anvi'o and pyani v0.2.11. Representative genomes were chosen based on quality scores and clustered based on >95% ANI. The final list of MAGs were taxonomically assigned with GTDB-Tk (Genome Taxonomy Database Toolkit; v2.1.0) using classify\_wf with default settings.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Research [guidelines for submitting code & software](#) for further information.

## Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A list of figures that have associated raw data
- A description of any restrictions on data availability

Human derived sample types were aligned to KneadData's default human reference genome (a combination of hg38 human genome reference (GenBank assembly accession # GCA\_000001405.29) and small contaminant sequences) and aligned reads were removed. Samples acquired from mice were processed using the reference genome GRCm39 constructed from C57BL/6J mouse-strains (GenBank assembly accession #GCA\_000001635.9). After bioinformatic host removal, the percentages of host reads were calculated by dividing reads remaining after host filtering by the total reads that passed QC. Metaphlan 3.0 workflow utilized database mpa\_v30\_CHOCOPhiAn\_201901.

For anvi'o workflow analysis, host reads were removed by alignment to the hg38 human genome reference (GenBank assembly accession # GCA\_000001405.29). The datasets generated during and analyzed during the current study will be made available upon publication at CaltechDATA, <https://doi.org/10.22002/gx69z-wec80>. Microbial sequencing data will be made available upon publication at NCBI Accession PRJNA991155.

The code utilized in data processing and analysis will be made available upon publication at CaltechDATA, <https://doi.org/10.22002/gx69z-wec80>.

## Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

Life sciences       Behavioural & social sciences       Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see [nature.com/documents/nr-reporting-summary-flat.pdf](https://nature.com/documents/nr-reporting-summary-flat.pdf)

## Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	No sample size calculation was performed. Minimal variations in MEM host removal across subjects (human and rodent) of any sample type were detected. Most analyses were performed on intraindividual controls (i.e. collection of multiple biopsies/saliva/scrapings/stool/etc. with and without MEM-treatment) to account for the small sample size.
Data exclusions	One human patient was excluded during analysis due to errors during the wet-lab processing with MEM. The MEM protocol was not performed correctly on this patient's samples, which was noted during the in-lab processing at University of Chicago. Exclusion criteria were not pre-established.
Replication	Triplicate biopsies were obtained from each region of the GI tract. Variation in technical replications of all measurements were shown to be smaller than the presented variation in biological replications. Analyses of technical (e.g., multiple replicates to determine highest possible PCR error rate) and biological replicates are reported in figure captions.
Randomization	Samples (rodent and human) were split randomly into control and MEM-treated groups.
Blinding	Investigators were not blinded during data collection and analyses

## Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

## Materials &amp; experimental systems

n/a	Included in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input checked="" type="checkbox"/>	<input type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input type="checkbox"/>	<input checked="" type="checkbox"/> Animals and other organisms
<input type="checkbox"/>	<input checked="" type="checkbox"/> Human research participants
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

## Methods

n/a	Included in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

## Animals and other organisms

Policy information about [studies involving animals](#); [ARRIVE guidelines](#) recommended for reporting animal research

Laboratory animals	Male and female wild-type, non-transgenic surplus mice were used in this study. One wild-type Syngap surplus rat was used for tissue collection. All animals were healthy prior to sacrifice. All animals were used exclusively for tissue collection, and intra-individual controls were used, therefore we do not expect housing conditions, light/day cycles, temperature, and humidity to affect the study's conclusions. These parameters were not collected.
Wild animals	The study did not involve wild animals.
Field-collected samples	The study did not involve samples collected in the field.
Ethics oversight	All animal husbandry and experiments were approved by the Caltech Institutional Animal Care and Use Committee (IACUC protocol #1769)

Note that full information on the approval of the study protocol must also be provided in the manuscript.

## Human research participants

Policy information about [studies involving human research participants](#)

Population characteristics	Adults scheduled for routine colon cancer screenings via colonoscopy at the University of Chicago Medicine (UCM) were screened for diagnosis and eligibility criteria for enrollment in the study on a weekly basis. Exclusion criteria included: participants with chronic infectious diseases such as human immunodeficiency virus (HIV) or hepatitis C (HCV); active, untreated <i>Clostridium difficile</i> infection; active infection with severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2); intravenous or illicit drug use such as cocaine, heroin, non-prescription methamphetamines; active use of blood thinners; severe comorbid diseases; participants on active cancer treatment; and participants who were pregnant. Samples arrived at Caltech de-identified with no information provided on age, gender, or genotypic information. We are unable to provide covariate analysis.
Recruitment	For human intestinal biopsies, approaching prospective participants was at the discretion of their treating physician and was not done in cases that would put participants at any increased risk, regardless of reason. Participants were approached the day of their procedure and informed, written consent was obtained before any samples were acquired. We do not expect any biases present in participant recruitment to impact the study's conclusions. For human saliva samples, prospective participants were recruited from the student body at Caltech. There may be bias for recruitment of young adults of high socioeconomic standing. As analysis of samples was limited to identifying the impact of the processing method (i.e. MEM treatment), we do not expect these biases to impact the study's conclusions.
Ethics oversight	All activities related to enrollment of participants, collection of samples, and sample analysis were approved by the University of Chicago Institutional Review Board (IRB) and performed under University of Chicago IRB protocols #15573A and #13-1080. De-identified samples from UC were received at Caltech and analyzed under Caltech IRB protocol #21-1083. Human saliva samples were acquired from two healthy adult volunteers and analyzed under California Institute of Technology Institutional Review Board (IRB) protocol #21-1092. All participants provided (digital) written informed consent prior to donation. No personal identifying information was collected at the time of consent and participant specimens were coded.

Note that full information on the approval of the study protocol must also be provided in the manuscript.