

1 **Convergent genomic responses of human gut bacteria to variations in industrialization**

2

3 **Authors & Affiliations**

4

5 Malte Rühlemann^{1,2,3} Lénárd L. Szánthó⁴, Silvio Waschina⁵, Lucas Moitinho-Silva^{1,2}, Laura K.
6 Mews², Joan F. Camarena², Hannah Jebens^{1,2}, John Costa^{1,2}, Vanessa Juimo^{1,6#}, Alain Fezeu^{1,6#},
7 Adwoa Agyei^{1,7#}, Mary Y. Afihene^{1,8#}, Shadrack O. Asibey^{1,9#}, Yaw A. Awuku^{1,10#}, Amoako
8 Duah^{1,11#}, Yvonne A. Nartey^{1,12#}, Fatimah Ibrahim^{1,13,14#}, Yvonne A. L. Lim^{1,13,15#}, Tan M. Pin^{1,13,14#},
9 Charles Onyekwere^{1,16#}, John Rusine^{1,17,18#}, Ivan E. Mwikarago^{1,18,19,20#}, John Baines^{21,22}, Andre
10 Franke², Gergely J Szöllösi^{4,23}, Ramnik Xavier^{24,25,26}, Eric J. Alm^{1,24,25,27§}, Mathieu
11 Groussin^{1,2,24,25,27§}, Mathilde Poyet^{1,21,24,25,27§}

12

13 ¹ Global Microbiome Conservancy, microbiomeconservancy.org

14 ² Institute of Clinical Molecular Biology, Kiel University & University Hospital Schleswig-Holstein,
15 Kiel, Germany

16 ³ Institute for Medical Microbiology and Hospital Epidemiology, Hannover Medical School,
17 Hannover, Germany

18 ⁴ Model-Based Evolutionary Genomics Unit, Okinawa Institute of Science and Technology
19 Graduate University, Okinawa, Japan

20 ⁵ Division of Food Technology, Institute of Human Nutrition and Food Science, Kiel University,
21 Kiel, Germany

22 ⁶ Institut de Recherche pour le Développement, Yaounde, Cameroon

23 ⁷ Department of Medicine and Therapeutics, University of Ghana Medical School and Korle Bu
24 Teaching Hospital, Accra, Ghana

25 ⁸ Department of Medicine, Kwame Nkrumah University of Science and Technology, Kumasi,
26 Ghana

27 ⁹ Catholic University College, Sunyani, Ghana

28 ¹⁰ Department of Internal Medicine and Therapeutics, School of Medical Sciences University of
29 Cape Coast, Cape Coast, Ghana

30 ¹¹ Department of internal medicine, University of Ghana Medical Centre, Legon, Accra, Ghana

31 ¹² Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm,
32 Sweden

33 ¹³ Centre for Innovation in Medical Engineering, University of Malaya, Kuala Lumpur, Malaysia

34 ¹⁴ Department of Molecular Medicine, Faculty of Medicine, Universiti Malaya, 50603, Kuala
35 Lumpur, Malaysia

36 ¹⁵ Department of Parasitology, Faculty of Medicine, Universiti Malaya, 50603 Kuala Lumpur,
37 Malaysia

38 ¹⁶ Department of Medicine, Lagos State University College of Medicine, Lagos, Nigeria

39 ¹⁷ National Institute of Allergy and Infectious Diseases (NIAID), Rockville, Maryland, United States

40 ¹⁸ National Reference Laboratory, Kigali, Rwanda

41 ¹⁹ Rwanda FDA and College of Medicine and Health Science of the University of Rwanda
42 affiliations and that of the Reference laboratory, Rwanda

43 ²⁰ Health Science of the University of Rwanda

44

45 ²¹ Institute of Experimental Medicine, Kiel University, Kiel, Germany

46 ²² Guest Group Evolutionary Medicine, Max Planck Institute for Evolutionary Biology, Plön,
47 Germany

48 ²³ HUN-REN Centre for Ecological Research, Institute of Evolution 1121 Budapest, Hungary

49 ²⁴ The Broad Institute of MIT and Harvard, Cambridge, MA, United States.

50 ²⁵ Center for Microbiome Informatics and Therapeutics, Massachusetts Institute of Technology,
51 Cambridge, MA, United States.

52 ²⁶ Center for Computational and Integrative Biology and Department of Molecular Biology,
53 Massachusetts General Hospital and Harvard Medical School, Boston, MA, United States.

54 ²⁷ Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA,
55 United States.

56

57 #These authors contributed equally to this work

58 §: co-senior authors

59 Corresponding authors: Mathilde Poyet, Mathieu Groussin and Eric Alm.

60

61

62 **Summary**

63

64 To what extent gut bacteria respond to the distinct ecological pressures imposed by human
65 lifestyle remains unclear. Here, we investigate how genomic adaptation in gut bacteria differ
66 between industrialized and non-industrialized human populations. We generated a broad
67 collection of isolate genomes spanning diverse host geographies, lifestyles, species, and strains.
68 We first found that compared to MAGs, paired isolate genomes recover more functional elements
69 and signals of horizontal gene transfers (HGTs). Leveraging isolate genomes from multiple
70 species, we find that strains from industrialized hosts experience an expansion of proteome size
71 and harbor greater pangenome fluidity, driven by recent events of HGTs. Gene- and variant-level
72 analyses reveal convergent patterns of lifestyle-specific adaptation in functions that are critical for
73 ecological adaptation, such as stress response, cell envelope remodeling and central metabolism.
74 Our results demonstrate that industrialization imprints evolutionary signatures on gut bacterial
75 genomes, illuminating the effects of rapidly changing environments on human biology.
76

77 Introduction

78

79 Human populations living in industrialized societies harbor gut bacterial communities that differ
80 markedly from those in non-industrialized populations¹⁻⁴. These differences are driven by
81 variation in lifestyle, diet, sanitation and exposure to environmental microbes^{5,6}. As such, lifestyle-
82 associated factors not only reshape the nutrient landscape available to gut bacteria, but also alter
83 the set of microbial interaction partners with which each species coexists and interacts. While
84 shifts in microbiome composition across human populations – particularly in relation to
85 industrialization and subsistence strategies – are now well documented, far less is known about
86 whether and how individual bacterial strains evolve and adapt to these host-associated
87 environments. In contrast to environmental microbes and pathogens⁷, adaptive genomic
88 signatures in commensal gut bacteria remain poorly characterized, despite evidence for
89 adaptation occurring within individual hosts on timescales of days to months⁸⁻¹².

90

91 Addressing this gap requires deep, high-quality collections of cultured isolates and their genomes
92 from a wide range of human lifestyles and geographic settings. However, existing culture
93 collections remain heavily biased, primarily representing microbiota from individuals living in
94 countries with high human development index (HDI) such as the United States, China, and
95 European nations¹³⁻²⁰. In our previous work²¹, we began to address this imbalance by generating
96 the first Global Microbiome Conservancy (GMbC) isolate genome collection, comprising over
97 4,000 gut bacterial genomes from host populations spanning a broad range of geographies and
98 lifestyles. Others also expanded the phylogenetic and geographic diversity of target taxa, such as
99 *Segatella copri* (previously named *Prevotella copri*)²².

100

101 In addition to isolates, metagenome-assembled genomes (MAGs) have become a widely used
102 resource to explore microbial diversity²³⁻²⁶, including from non-industrialized populations for
103 which access to isolate genomes is more challenging. While MAGs have enabled large-scale
104 surveys of phylogenetic diversity, their use for detecting recent adaptation or fine-scale genomic
105 features remains debated^{27,28}. Potential limitations of MAGs include variable completeness, the
106 inability to fully capture within-host population heterogeneity, and the loss of accessory functions
107 – particularly those associated with mobile genetic elements (MGEs) and horizontal gene
108 transfers (HGTs)^{27,29}. These challenges stem from technical and biological factors: binning errors,
109 reliance on single-copy core genes (SCGs) for estimating completeness, and the inherent
110 difficulty of assembling genomes in the presence of high microdiversity and strain-level variation
111^{30,31}. To robustly evaluate MAG quality, direct comparison of MAGs to taxonomically-paired isolate
112 genomes from the same sample is required. While efforts have been made in this direction²⁸,
113 such comparisons have so far been limited in scale and taxonomic breadth, leaving a critical gap
114 in our ability to assess how well MAGs represent real genomic diversity.

115

116 Here, we expand and leverage the GMbC collection of isolate genomes and MAGs to investigate
117 how industrialization shapes the evolutionary trajectories and adaptive potential of gut commensal
118 bacteria. We found that gut bacterial genomes from industrialized hosts show evidence of
119 proteome expansion and elevated pangenome fluidity, both owing to a recent acceleration in
120 HGT-driven gene acquisition. Across multiple species, we uncover parallel signals of genomic

121 adaptation to host industrialization status, including lifestyle-specific gene enrichment, signatures
122 of positive selection, and convergent non-synonymous single nucleotide variants (SNVs) in genes
123 that are functionally relevant for ecological adaptation.

124 Results

125

126 Extensive species, strain and geographic diversity in the GMbC isolate genome collection

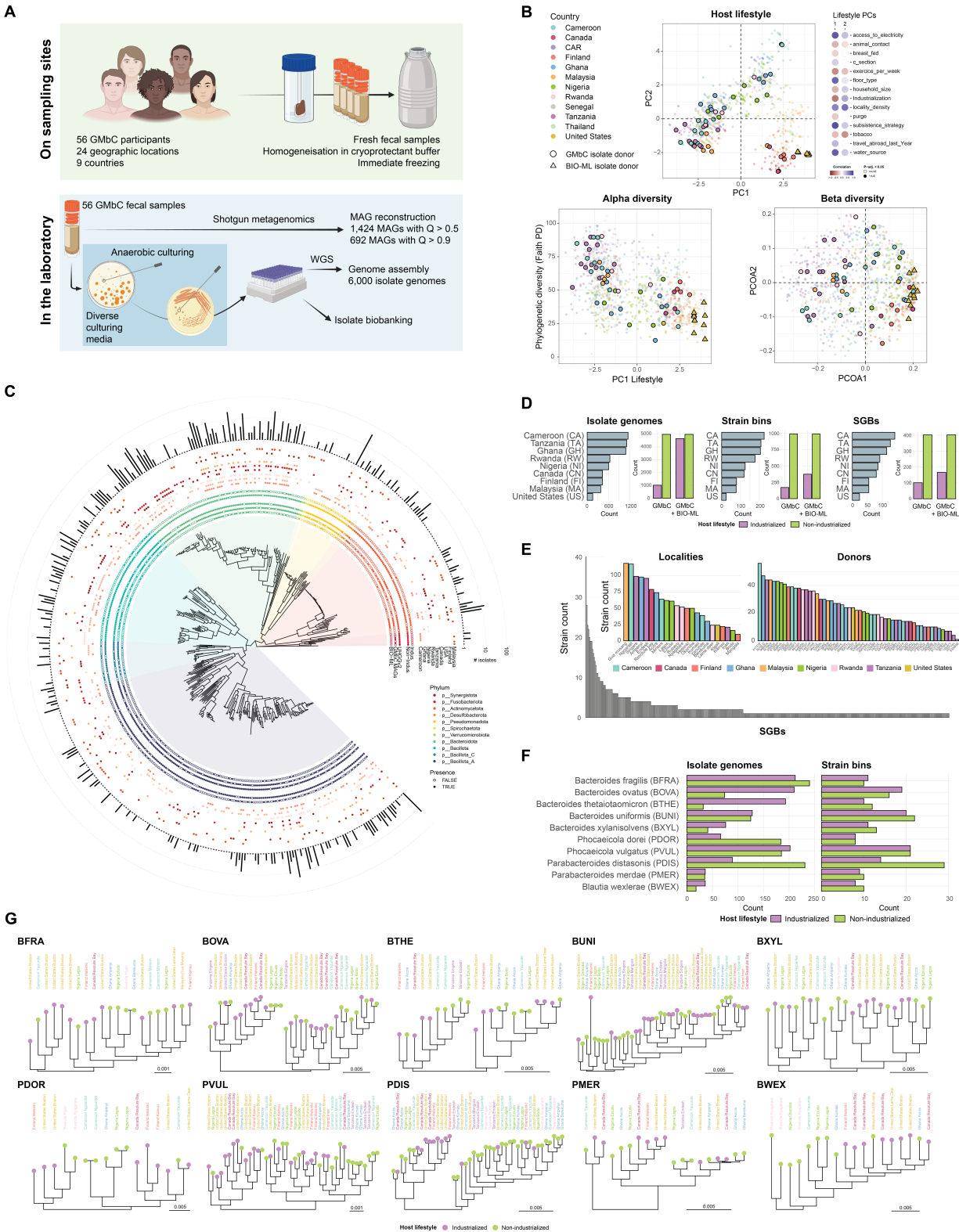
127

128 To build the GMbC isolate genome collection of human gut bacteria, we employed culturing
129 strategies designed to maximize bacterial species and strain diversity (see Methods) (Fig. 1A).
130 Host individuals from the GMbC cohort ⁶ were selected to represent a broad range of lifestyles,
131 geographic backgrounds and microbiome diversity (see Methods) (Fig. 1B & Supp. Table 1). We
132 release a new set of 1,841 high-quality isolate genomes, which we integrate into our previously
133 published set of 4,140 genomes ²¹, resulting in a total of 5,981 high-quality genomes (see
134 Methods) (Fig. 1C). GMbC genomes have a median completeness of 99.2%, contamination of
135 1.67%, and quality score of 98.3. Using similarity thresholds at 99% and 95%, GMbC isolate
136 genomes were clustered into 1,133 strain- and 434 species-level genome bins (StGBs and SGBs,
137 respectively). Strain-level genome bins represented in multiple hosts were further split by host to
138 ensure that each StGB represented a host-specific strain (see Methods). These isolates were
139 cultured from stool samples of 56 donors across 24 geographic locations in 9 countries (Fig. 1C,
140 Supp. Fig. 1, and Supp. Table 1). Of the 434 SGBs, 403 were obtained from individuals living
141 more non-industrialized lifestyles, and 102 from individuals in more industrialized settings (Fig.
142 1D). For comparison, only 106 SGBs were recovered in our previous BIO-ML collection, which
143 included only individuals from industrialized populations in the United States ¹². The median
144 number of GMbC SGBs per donor, geographic location, self-declared ethnicity, and country is 17,
145 41, 52, and 87, respectively (Fig. 1D & E).

146

147 Next, we assessed the phylogenetic and functional diversity of the GMbC isolate genome
148 collection. The 434 GMbC SGBs span 13 distinct phyla and 46 families, encompassing both
149 dominant taxa of the human gut microbiome and phylogenetic groups that are rare,
150 underrepresented, or difficult to cultivate (Fig. 1C). Notably, the collection includes isolates from
151 Spirochaetota (*Treponema_D succinifaciens*, n = 5; *Treponema_D peruense*, n = 2) and
152 Synergistota (*Cloacibacillus porcorum*, n = 6). Contrarily to *Treponema pallidum*, which is a
153 pathogen, very little is known about *Treponema succinifaciens*. It was frequently found to be in
154 relatively high abundance in the gut microbiome of non-industrialized populations, while being
155 absent from those of industrialized populations ^{6,32,33}. *Treponema* in non-industrialized
156 populations is thought to be a commensal that may promote fiber degradation in the context of
157 fiber-rich diets ¹. Looking at the broader GMbC cohort, we recently found *Treponema*
158 *succinifaciens* to be negatively associated with intestinal inflammation markers ⁶, and to be
159 negatively associated with hypoxia and TNF signaling gene expression in industrialized contexts
160 ³⁴. The *C. porcorum* isolates represent the first cultured strains of human origin, as previously
161 available reference strains (e.g., those in the DSMZ collection) were isolated from pigs ³⁵.
162 Furthermore, human-derived *C. porcorum* genomes are absent from large-scale metagenomic
163 resources such as the UHGGv2 MAG collection ^{24,36} (Fig. 1C). The GMbC dataset also includes
164 archaeal isolates, with 7 genomes spanning the phyla Methanobacteriota (*Methanobrevibacter_A*
165 *smithii*, n = 2; *Methanosphaera stadtmanae*, n = 1) and Thermoplasmata (*Methanarcanum*
166 *hacksteinii*, n = 4). Among the 434 SGBs, 55 lack a species-level taxonomic assignment,

167 highlighting gaps in current reference databases such as GTDB (rel214)³⁷. Of these, 27 belong
 168 to the *Collinsella* genus, 2 to *Prevotella*, and 2 to *Faecalibacterium*.
 169



170

171 **Figure 1 – Geographic, species, strain, and host lifestyle diversity in the GMbC collection**
172 **of human gut bacterial isolate genomes**

- 173 A. Overview of sampling, preservation, culturing, isolation, and sequencing procedures for
174 gut bacterial genomes (see Methods).
- 175 B. Lifestyle and microbiome diversity of donors used for culturing and isolating gut bacteria
176 in the context of the broader GMbC + BIO-ML cohort. Top panel: dimensional reduction
177 analysis of various lifestyle factors (see Methods). Donors used for culturing are shown in
178 larger symbols with dark border. GMbC donors are shown in circles, BIO-ML donors are
179 shown in triangles. Spearman correlations between the first two PCs and individual
180 lifestyle factors are shown on the right. Alpha diversity (measured with Faith PD index)
181 and beta diversity (unweighted UniFrac) of GMbC and BIO-ML isolate donors and
182 participants are shown in bottom panels.
- 183 C. Phylogenomic tree of representative genomes from 434 species-level genome bins
184 (SGBs). Inner ring shows overlap with external genome collections (UHGG v2, GMbC
185 MAGs, BIO-ML). Middle ring indicates host lifestyle origin (industrialized or non-
186 industrialized). Outer ring shows country distribution and isolate genome counts per SGB.
187 Clade colors represent phyla.
- 188 D. Isolate genome, strain bin, and SGB counts by country and host lifestyle. Strain bins group
189 genomes from the same donor with >99% similarity (see Methods). Counts that include
190 isolate genomes of the BIO-ML collection per host lifestyle are also shown. BIO-ML isolate
191 genomes were generated following the same pipeline as described in A (see Methods).
- 192 E. Distribution of strain bin counts across SGBs, localities and individual hosts. Colors denote
193 country.
- 194 F. Ten bacterial species with ≥ 8 strain bins sampled from industrialized or non-industrialized
195 hosts. Barplots show isolate and strain bin counts per lifestyle.
- 196 G. Phylogenetic trees of representative strain bin genomes for the 10 species in panel E. Tip
197 points indicate host lifestyle; labels show country/locality and are color-coded by country.
198 Trees are midpoint-rooted. Branch length scales are in expected number of substitutions
199 per site.

200

201 The GMbC collection also captures substantial within-species strain genomic diversity (Fig. 1D &
202 E & Supp. Fig. 1). In total, 39 SGBs are represented by at least 10 distinct StGBs, including from
203 multiple key bacterial species in the human gut, such as *Bacteroides*, *Parabacteroides*, *Blautia*
204 species. Genomic-based phenotypic and functional predictions, including amino acid
205 auxotrophies that are associated with chronic diseases³⁸, suggest substantial strain-level
206 diversity among several key taxa, including *Prevotella*, *Veillonella* and *Blautia* species (Supp. Fig.
207 1). Overall, the GMbC collection extensively samples species and strain-level diversity of human
208 gut bacteria across various human host modalities, including geography and lifestyle, providing
209 unprecedented amounts of genomic material for in-depth genomic and functional investigations
210 of the global gut microbiome. We leverage this phylogenomic diversity in the following analyses
211 to investigate how host industrialization influences the genomic evolution of human gut bacteria.

212
213

214 **Isolate genomes recover more functional and mobile elements than MAGs**

215

216 We first benchmarked our isolate genomes against metagenome-assembled genomes (MAGs)
217 to assess the advantage of combining MAGs with isolate genomes in our analysis. Previous
218 studies have questioned the completeness and quality of MAGs due to errors introduced during
219 metagenomic assembly and binning, as well as challenges such as within-sample strain diversity
220 ³⁹. Although adding MAGs could substantially expand strain representation and species coverage,
221 we wondered whether these hypothesized drawbacks of MAGs could impact our analysis of
222 genomic evolution across host lifestyles, particularly for accessory gene families and mobile
223 genetic elements that may be central to lifestyle-driven adaptive processes. So far, it has
224 remained difficult to systematically evaluate how MAGs compare to isolate genomes because
225 most prior comparisons involved unpaired genomes—i.e., MAGs and isolates obtained from
226 different hosts, and often representing different strains. An advantage of our study design is the
227 ability to directly compare the characteristics of isolate genomes with those of paired MAGs. To
228 do this, we leveraged shotgun metagenomic data that we recently generated from the same fecal
229 samples used to culture our isolates ⁶, and reconstructed MAGs using a multi-binning strategy
230 ^{6,40}. We filtered out low-quality MAGs using similar thresholds for completeness (<50%) and
231 contamination (>10%) as for our isolate genomes. This allowed us to assemble a unique dataset
232 of 147 paired MAG–isolate genomes, originating from the same species, and sampled from the
233 same donor (Fig. 2A). These pairs span a broad range of bacterial taxonomies and abundances
234 in the human gut (Fig. 2B & Supp. Fig. 2).

235

236 The median assembly quality score of the 147 MAGs is 0.90, comparable to the UHGGv2 MAG
237 collection (MAG_Q = 0.89) (see Methods). The MAGs show a median completeness of 92.5%
238 and contamination of 2.5% (Supp. Table 2). Among high-quality MAGs (MAG_Q > 0.9; n = 74),
239 the median MAG_Q reaches 0.96 (median completeness = 98% and contamination = 2.5%) –
240 close to that of their paired isolate genomes (Q = 0.98) (Fig. 2A). We used genome-scale
241 metabolic models (GMMs) and reference databases to profile all genomes for functional
242 categories and mobile genetic elements (MGEs) (see Methods).

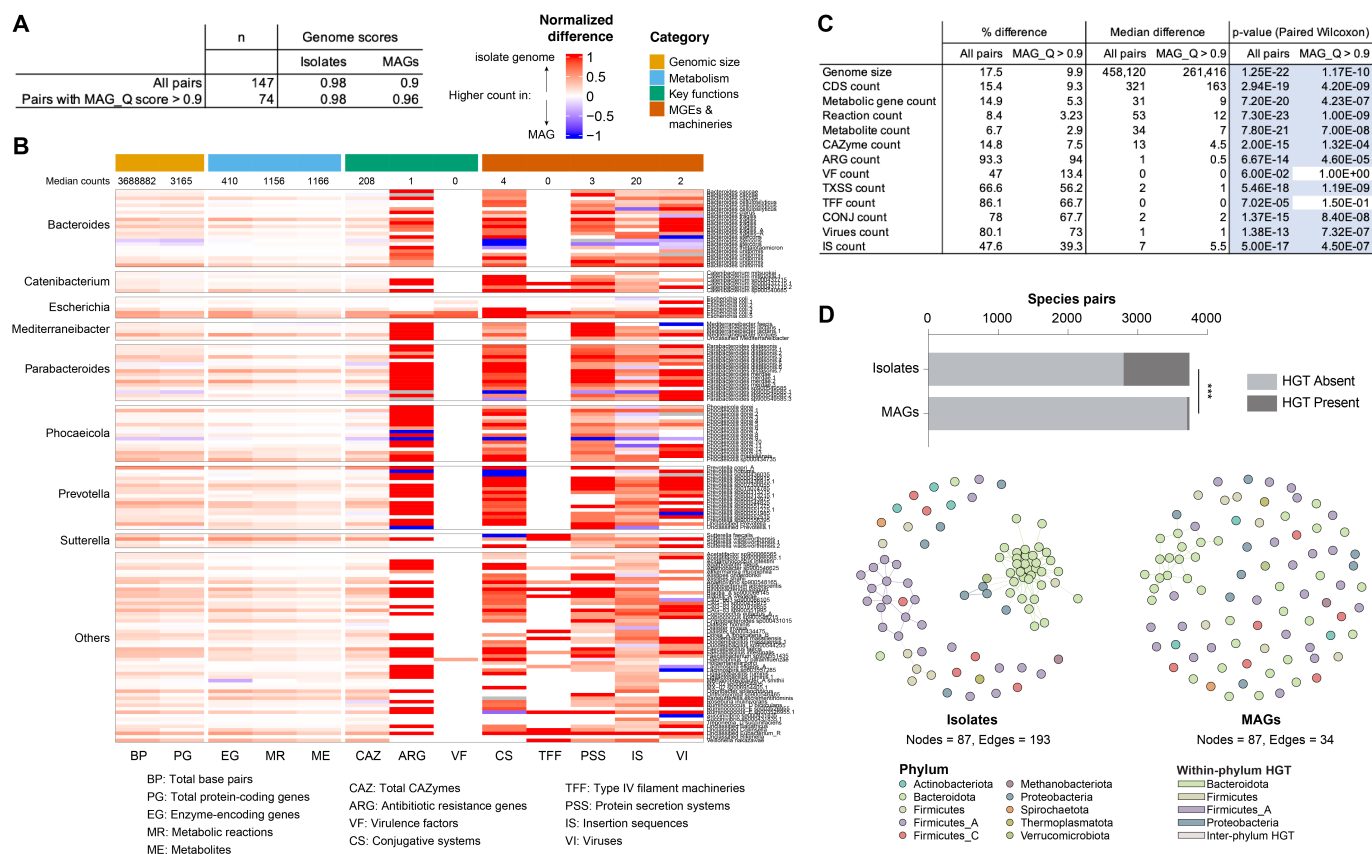
243

244 We found that isolate genomes are significantly larger (paired Wilcoxon test, median difference =
245 458,120bp, p = 3.6e-22) and contain more coding sequences (CDSs) (median diff. = 321 genes,
246 p = 2.9e-19) than MAGs (Fig. 2B & C, Supp. Fig. 2 & Supp. Table 2). Isolate genomes also harbor
247 more metabolic features predicted by GMMs, including enzyme-encoding genes (median diff. =
248 31, p = 7.2e-20), reaction pathways (median diff. = 53, p = 7.3e-23), and predicted metabolites
249 (median diff. = 53, p = 7.3e-23) (Fig. 2B & C). This disparity is further pronounced when GMM
250 reconstruction is performed without pathway gap-filling (see Supp. Table 2 and Methods). We
251 also found that MAGs required significantly more pathway gap-filling during GMM reconstruction
252 than isolate genomes (p = 4.7e-12), confirming that gene content inference and metabolic
253 reconstructions are more fragmented in MAGs.

254

255 We then compared the counts of specific gene families involved in key bacterial functions, such
256 as carbohydrate metabolism (carbohydrate active enzyme [CAZyme]), antibiotic resistance
257 (antibiotic resistance genes [ARGs]), and virulence (virulence factors [VFs]). These gene families

258 were consistently detected in greater numbers in isolate genomes compared to MAGs ($p = 2.0e-$
 259 15 , $p = 6.7e-14$ and $p = 6e-02$, respectively) (Fig. 2B & C, Supp. Fig. 2).
 260



261
 262 **Figure 2 – Isolate genomes recover more genomic features and HGT events than MAGs.**

263 A. Number and quality scores of MAG–isolate genome pairs. Pairs originate from the same
 264 donor sample and species.
 265 B. Heatmap comparing genomic feature counts across all genome pairs. Genera and species
 266 of genome pairs are shown on the left and right side of the heatmap, respectively. Genomic
 267 features are shown in columns, and are grouped in four categories: genomic size,
 268 metabolism, key functions and mobile genetic elements (MGEs) & machineries. For each
 269 pair, the difference in counts between the isolate genome and the MAG was calculated
 270 and normalized to the count in the isolate genome. Features with higher counts in the
 271 isolate genome or in the MAG are shown along a gradient of red to blue, respectively.
 272 C. Summary statistics of feature differences across all pairs.
 273 D. Comparison of HGT events. Between-species horizontal gene transfers (HGTs) were
 274 detected across isolate genomes, and across MAGs separately (see Methods). Genomes
 275 of MAG-isolate genome pairs cluster in 87 SGBs. Ratio of species pairs with detected
 276 HGTs ($n \geq 1$ HGT) were compared with a proportion test (Two proportion Z-test, ***: $p =$
 277 $7.96e-33$). Edges in the network indicate that at least 1 HGT was detected between
 278 species (nodes).
 279

280 We also profiled mobile genetic elements and machinery involved in horizontal gene transfer,
281 including conjugative elements (Conj), insertion sequences (ISs), prophages, type IV filament
282 (TFF) superfamily proteins, and protein secretion systems (TxSS). Each of these categories was
283 more abundant in isolate genomes compared to MAGs (Fig. 2B & C, Supp. Fig. 2).

284

285 Importantly, these trends remained consistent even when restricting the analysis to high-quality
286 MAGs (MAG_Q > 0.9; n = 74) (Fig. 2C, Supp. Table 2). Finally, differences in functional content
287 between isolates and MAGs were robust across a range of bacterial taxonomies and relative
288 abundances (Fig. 2B, linear models with taxonomy and abundance, see Supp. Table 2).

289

290 Considering that MGEs are better captured in isolate genomes, we next asked whether isolates
291 would also provide greater sensitivity and accuracy for detecting recent horizontal gene transfer
292 (HGT) events. To test this, we applied a previously established BLAST-based pipeline for
293 identifying nearly identical sequences shared between genomes of different species^{21,41}. Our 147
294 MAG–isolate pairs represent 87 distinct bacterial species. We used this pipeline to detect HGT
295 events across these 87 species using either MAGs or isolate genomes (Methods). Both the
296 number of candidate HGTs (BLAST hits) and the proportion of species pairs involved in putative
297 HGT were significantly higher when using isolate genomes (paired Wilcoxon test, $p = 3.27e-39$;
298 Two proportion Z-test, $X\text{-squared} = 142.4$, $p = 7.96e-33$) (Fig. 2D).

299

300 Overall, isolate genomes consistently recover higher numbers of genomic and functional features
301 compared to MAGs. Our results also show that high-quality MAGs still capture much of this
302 information (Fig. 2B & Supp. Table 2). With ultra-deep sequencing, long-read metagenomics and
303 advances in binning algorithms^{2,42,43}, MAGs should achieve performance comparable to isolates
304 in detecting these features. In the following, we chose to consider MAGs for validation analyses:
305 all primary inferences regarding gene content variation and sequence divergence are drawn from
306 the more complete isolate genomes, and MAGs are used when needed to confirm that these
307 patterns persist when sampling a broader set of host individuals (Fig. 7). By leveraging the
308 complementary strengths of both data types, we aim to minimize the risk that assembly artifacts
309 bias our conclusions about lifestyle-associated genomic evolution.

310

311

312 **Bacterial SGBs with broad sampling of StGBs to study the effect of host industrialization**
313 **status on genomic evolution**

314
315 To test whether host industrialized vs. non-industrialized lifestyles exert selective pressures that
316 impact genomic evolution and adaptation of gut bacteria, we investigated a range of genomic and
317 phylogenetic features, including gene content, ancestral gene gain and loss events, signals of
318 positive selection, and individual SNVs. We focused on ten bacterial species that were well-
319 represented in the GMbC collection and included a high number of isolate genomes and StGBs
320 from individuals of both lifestyle types (Fig. 1F & G): *Bacteroides fragilis* (BFRA), *Bacteroides*
321 *ovatus* (BOVA), *Bacteroides thetaiotaomicron* (BTHE), *Bacteroides uniformis* (BUNI),
322 *Bacteroides xylanisolvens* (BXYL), *Phocaeicola dorei* (PDOR), *Phocaeicola vulgatus* (PVUL),
323 *Parabacteroides distasonis* (PDIS), *Parabacteroides merdae* (PMER), and *Blautia A wexlerae*
324 (BWEX). We complemented the GMbC genomes with published isolate genomes of these
325 species from the BIO-ML collection (industrialized host donors from the Boston (MA, USA) area)
326 that we generated using similar culturing and sequencing methods¹² (Fig. 1F). For each species,
327 we reconstructed the pangenome (core, accessory and cloud genomes) and the recombination-
328 aware phylogeny of StGB representative genomes (Fig. 1G and Supp. Fig. 3) (see Methods).
329 Phylogenetic signal analyses of the industrialization trait (industrialized / non-industrialized host)
330 using Blomberg's K revealed that the trait has limited signal in most species ($K < 1$; Supp. Table
331 1) and is broadly distributed across StGB phylogenies (Fig. 1G). Only BWEX showed a significant
332 phylogenetic signal ($K > 1$, $p\text{-val} < 0.01$). Combined with appropriate control for phylogenetic
333 structure in statistical models, these patterns enable the detection of convergent genomic
334 responses to host industrialization status among distantly-related strains, as explored in the
335 following sections.

336
337 While we use a binary variable for industrialization status (industrialized vs. non-industrialized
338 host) to investigate patterns of bacterial adaptation, **we acknowledge that this classification**
339 **oversimplifies a broad spectrum of lifestyle differences**. In particular, non-industrialized
340 populations exhibit greater diversity in subsistence strategies and environmental exposures
341 compared to industrialized populations⁶. **To address this limitation, we also incorporate a**
342 **continuous variable derived from a multidimensional reduction of lifestyle-related factors**
343 **(named “PC1 Lifestyle” thereafter)**, including but not limited to industrialization status (Fig. 1B).
344 This variable was computed from the broader GMbC cohort ($n = 1,015$ participants) that we
345 recently described⁶, which includes the GMbC participants used here for bacterial isolation. PC1
346 Lifestyle is strongly correlated with industrialization status, while also capturing finer-scale
347 variation in lifestyle across participants (Fig. 1B). Where appropriate (i.e. gene enrichment and
348 single nucleotide variant (SNV) analyses in Fig. 5 & 7), we use this higher-resolution variable as
349 a continuous proxy for industrialization to validate associations identified using the binary
350 classification.

351
352
353

354 **Proteome expansion and increased pangenome fluidity in gut bacteria of hosts with**
355 **industrialized lifestyles**

356

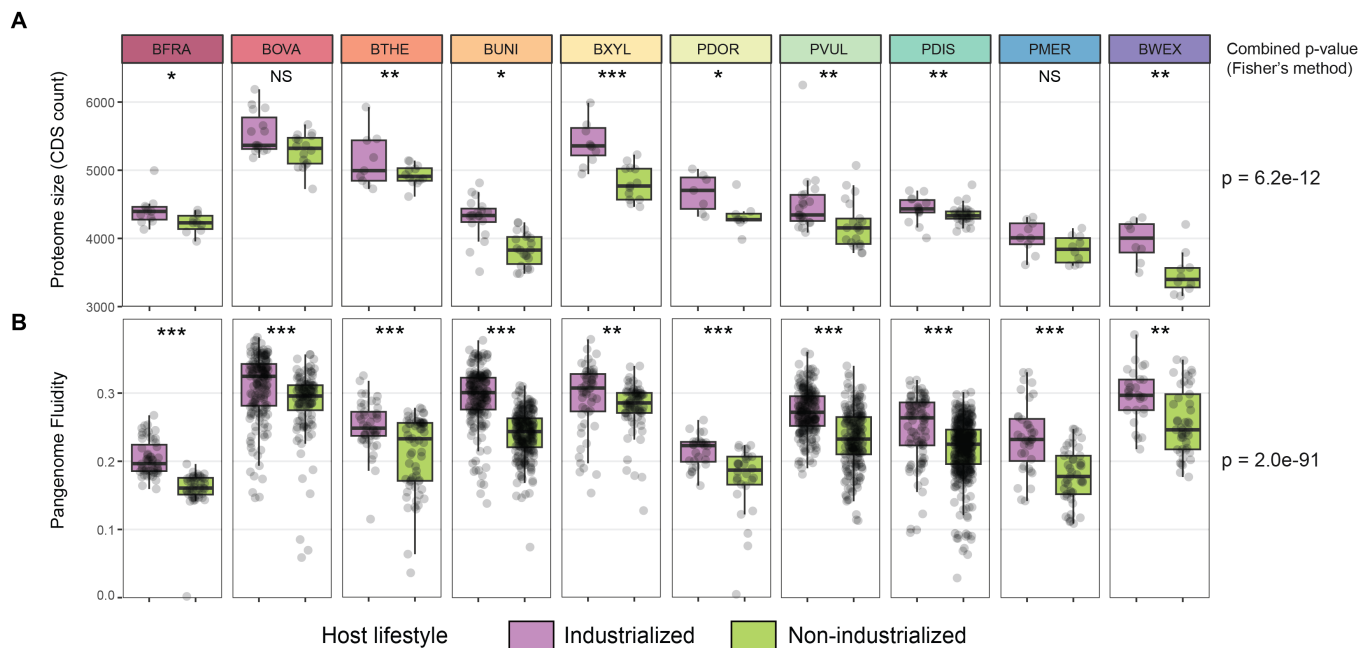
357 We first tested whether the total protein-coding gene content differs between strains colonizing
358 individuals from industrialized vs. non-industrialized lifestyles. Across all 10 species examined,
359 strains from industrialized hosts consistently exhibited larger proteome sizes than those from non-
360 industrialized hosts (Fig. 3A). To account for phylogenetic structure, we applied phylogenetic
361 linear regression models across species and identified eight species with significantly larger
362 proteomes in industrialized strains ($p < 0.05$). Combining species-level p-values using Fisher's
363 method revealed a strong overall signal of proteome expansion in industrialized populations ($p =$
364 $6.2e-12$) (Fig. 3A).

365

366 This proteome expansion may be associated with increases in pangenome size and fluidity. We
367 first confirmed that core genome size does not differ between strains from industrialized and non-
368 industrialized hosts (paired Wilcoxon test, $p = 0.92$), as expected, indicating that the observed
369 proteome expansion is not driven by changes in conserved genes. We then assessed pangenome
370 fluidity, which quantifies the variability in gene content among strains of the same species.
371 Specifically, fluidity measures the proportion of genes that are not shared between genome pairs,
372 relative to the total number of genes (see Methods). Higher pangenome fluidity reflects a more
373 open and dynamic pangenome, characterized by extensive gene turnover and accessory genome
374 diversity. Across all ten species analyzed, we consistently observed higher pangenome fluidity
375 among strains from industrialized hosts ($p\text{-val} < 0.01$ for each species; Fisher's combined $p\text{-val} =$
376 $2.0e-91$) (Fig. 3B), suggesting that bacterial strains in industrialized environments experience
377 stronger gene turnover and greater genomic plasticity. These findings align with our previous work
378 showing elevated rates of horizontal gene transfer (HGT) in gut bacteria from industrialized
379 populations²¹.

380

381



382

383

384 **Figure 3 – Industrialized host strains exhibit larger proteomes and signatures of relaxed**
 385 **selection**

386 A. Comparison of proteome size (coding gene counts) between strains of host with
 387 industrialized vs. non-industrialized lifestyles (in purple and green, respectively) across
 388 the 10 species presented in Figure 1. Counts were statistically compared while accounting
 389 for phylogeny (phyloglm function, see Methods) (***: p-value < 0.001; **: p-value < 0.01;
 390 *: p-value < 0.05; NS: non-significant – this legend applies to all other panels). P-values
 391 were combined with the Fisher’s method to test for cross-species evidence of differences
 392 in proteome size against the null hypothesis. This p-value is shown on the right of the
 393 panel.

394 B. Comparison of pangenome fluidity among industrialization- and non-industrialization-
 395 associated strains. The ratio of shared genes was calculated for strain bin pairs, using
 396 representative genomes. P-values were combined with the Fisher’s method (p-value
 397 shown on the right of the panel)
 398

399 **Recent increase in rates of bacterial gene gains among industrialized hosts**

400

401 We next sought to test whether the observed increase in proteome size among strains from
402 industrialized hosts is driven by elevated rates of gene acquisition or by reduced levels of gene
403 loss, relative to strains from non-industrialized hosts. Moreover, it remains unclear whether this
404 expansion reflects a recent evolutionary response to industrialized lifestyles. If the latter were the
405 case, we would expect gene gains and proteome size increases to be confined to terminal
406 branches of bacterial species trees, corresponding to more recent evolutionary events.
407 Alternatively, these events could have occurred along ancestral lineages, with ecological niche
408 selection favoring the colonization of industrialized hosts by strains of larger proteomes.

409

410 To distinguish between these scenarios, we used a species tree-gene tree phylogenetic
411 reconciliation method implemented in AleRax⁴⁴ to identify the branches of the species trees along
412 which genes were gained or lost. AleRax distinguishes two processes of gene gain, either via
413 horizontal transfer or by origination. In order to reconstruct gene trees, we considered all gene
414 families being present in at least 4 different genomes across StGBs. We inferred whether
415 ancestral nodes along the species phylogeny were from industrialized or non-industrialized hosts
416 using Wagner parsimony. We then quantified gene gain/loss events and copy numbers across
417 internal and terminal branches based on reconciliation scenarios, stratified by host lifestyle (See
418 Methods, Supp. Table 4). We found higher counts of gene gains along branches associated with
419 industrialized lifestyles, that gene gains outnumbered losses across most species, and that they
420 mainly occurred along terminal branches (Fig. 4A, Wilcoxon-tests, plain circles show species with
421 statistically significant differences ($p\text{-val} < 0.05$) in event counts between lifestyles). Interestingly,
422 the higher rate of gene gain in *B. thetaiotaomicron* is amplified by high rates of gene gains along
423 ancestral branches associated with industrialized lifestyle, and higher rates of gene loss along
424 ancestral branches associated with non-industrialized lifestyle (Fig. 4A).

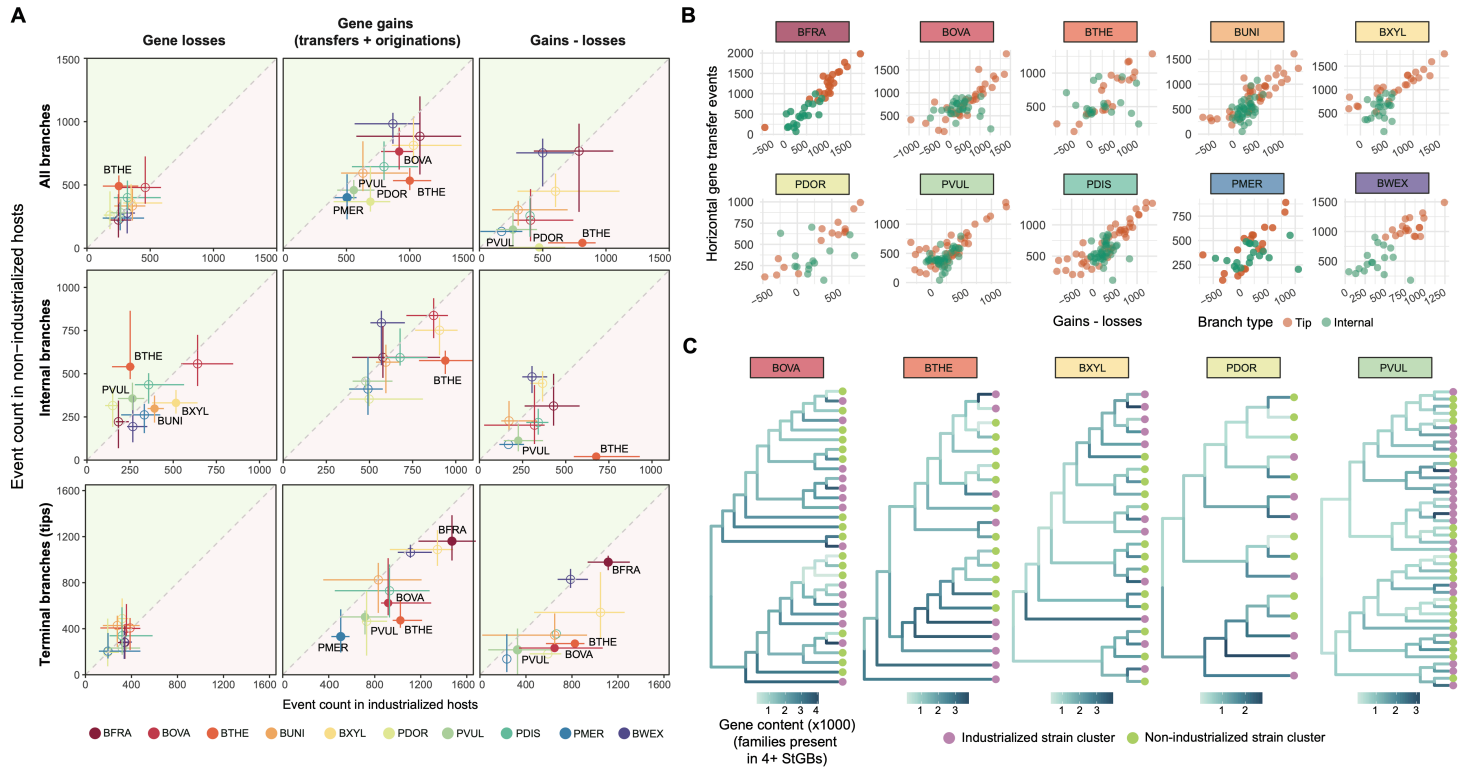
425

426 We further found that the difference between gene gain and loss counts across all branches is
427 strongly correlated with the number of HGT events per branch, supporting HGT as the primary
428 driver of gene acquisition in these lineages (Fig. 4B). Additionally, reconciliation-based
429 reconstructions of gene copy number along the species tree revealed that increases in gene
430 content predominantly occur along terminal branches associated with industrialized hosts (Fig.
431 4C and Supp. Fig. 4). Notably, gene families consisting of singletons or occurring in less than 4
432 StGBs were excluded from this analysis, indicating that proteome expansion is not solely driven
433 by the acquisition of rare genes, but also involves genes that are more broadly shared across
434 strains.

435

436 Altogether, these results show that the expansion of bacterial proteomes in industrialized hosts is
437 a relatively recent phenomenon, likely associated with the emergence of industrialized lifestyles,
438 and is mainly driven by high occurrence of HGTs.

439



440

441 **Figure 4 – Recent and HGT-driven gene gains promote proteome expansion in**
 442 **industrialized strains**

443 A. Species tree - gene tree reconciliations were sampled to detect and count per-branch
 444 events of gene transfer, loss, origination and speciation (see Methods). Counts of per-
 445 branch gene loss (left column) and gain (middle), and differences between gain and loss
 446 counts (right) were compared between host lifestyle categories (industrialized: purple
 447 area; non-industrialized: green area). Gene gains were defined as the sum of gene
 448 transfer and origination events. Top row: counts aggregated across all branches. Middle
 449 row: Counts of internal branches. Bottom row: counts of terminal (tip) branches. For each
 450 species, median counts are shown, with intervals ranging from the 25th to the 75th
 451 quantiles. Plain points indicate species for which the difference in counts between host
 452 lifestyle categories is significantly different (Wilcoxon tests). Species are colored-coded.

453 B. Correlation between per-branch gain–loss differences and HGT counts, broken down by
 454 internal (green) and terminal (orange) branches. All correlations are statistically significant
 455 (p -val < 0.001; Spearman correlation tests).

456 C. Increasing gene content along lineages of industrialized hosts. The panel depicts the
 457 evolution of the number of genes along the phylogeny of BOVA, BTHE, BXYL, PDOR and
 458 PVUL, based on the reconciliation-aware reconstruction of ancestral gene contents.
 459 Reconciliations were calculated from the set of gene families present in at least four StGBs
 460 (tips of the tree). These 5 species have significant differences in proteome size (Figure 3)
 461 and in gene gains along terminal branches between host lifestyles. Data for the other five
 462 species, which show similar trends, is presented in Supp. Fig. 4.

463

464 **Differential enrichment of genes based on host industrialization status**

465

466 Building on our finding that host industrialization influences pangenome size and fluidity, we
467 hypothesized that specific gene families within the accessory genome may exhibit differential
468 prevalence between strains from industrialized and non-industrialized hosts. To test this, we
469 analyzed gene presence/absence patterns across the pangenome of each species (mean = 2,691
470 gene families; SD = 1,515), using phylogeny-aware logistic regression models and the
471 recombination-aware StGB phylogenies. We considered host industrialization status as a binary
472 response variable, and significant associations were further validated using PC1 Lifestyle as a
473 continuous response variable (see Fig. 1B and Methods) for host industrialization. On average,
474 we found that 5.13% (SD = 8.32%) of gene families were significantly associated with host lifestyle
475 (FDR-adjusted $p < 0.05$) (Fig. 5A, Supp. Fig. 5 & Supp Table 5).

476

477 To assess whether gene-level associations reflect convergent adaptation across species, we
478 identified genes with significant industrialization-associated prevalence in multiple taxa. Such
479 convergent signals, replicated both within and between species, would support the hypothesis of
480 positive selection acting on specific genes in response to lifestyle-related environmental
481 pressures. We found six gene families with significant associations in two or more species, with
482 all but one gene showing consistent direction of enrichment direction in relation to
483 industrialization: *cas1* (BFRA, PDIS), *int* (BUNI, PDIS), *ugd* (BUNI, PDIS), *wecA* (BOVA, PDIS),
484 *per1* (BUNI, PDIS, PVUL), and *traM* (BOVA, BUNI, PDIS) (Fig. 5B). Interestingly, both *ugd* and
485 *wecA* are involved in the formation of bacterial cell surface structures, particularly in the synthesis
486 of polysaccharides and glycans that form key components of the envelope, such as capsule, O-
487 antigen LPS, and enterobacterial common antigens^{45,46}.

488

489 Expanding to functional annotations, we identified 23 KEGG Ortholog (KO) groups with significant
490 prevalence associations in at least two species. Of these, four appeared in three or more species:
491 K01185 (lysozyme; BOVA, BUNI, PDIS); K07154 (serine/threonine-protein kinase HipA; BFRA,
492 PDIS, PVUL); K17836 (beta-lactamase class A; BUNI, PDIS, PVUL); and K21572 (starch-binding
493 outer membrane protein SusD/RagB; BFRA, BOVA, BUNI, BXYL, PDIS, PDOR) (Supp. Table 5).

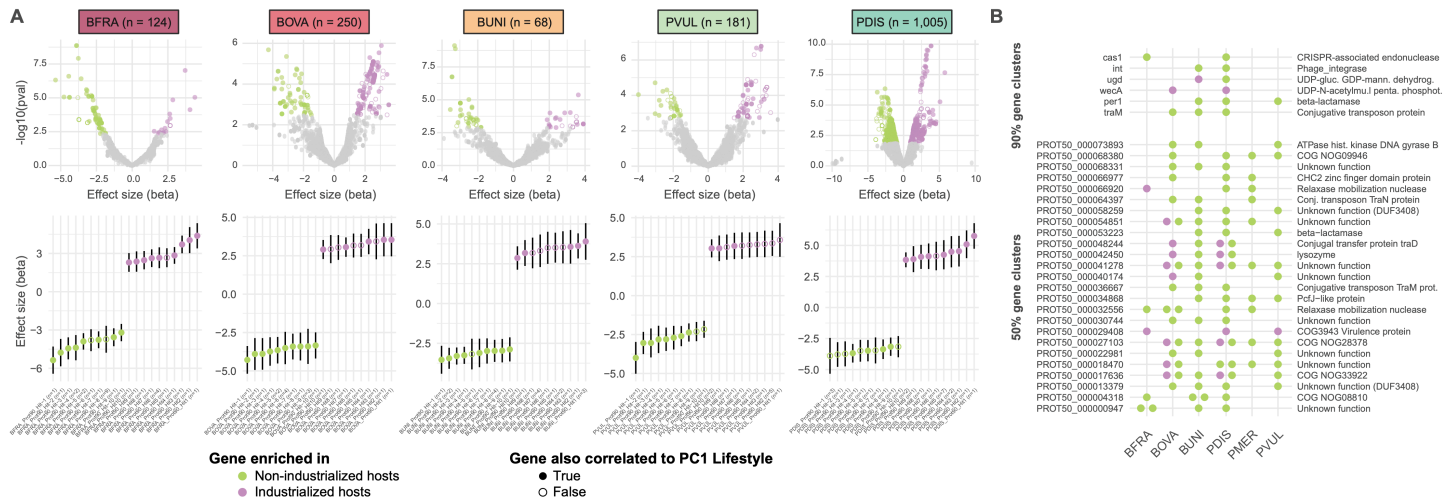
494

495 To further explore cross-species convergence, we clustered the species-level catalogs of 90%-
496 similarity gene families into a broader, cross-species catalog at 50% sequence similarity. This
497 revealed 25 gene families with industrialization-associated signals in multiple species (Fig. 5B &
498 Supp. Table 5). Although many of these lacked annotations, we found that four belong to the *tra*
499 operon, involved in conjugative transfer of genetic material (*traD*, *traK*, *traM*, *traN*)⁴⁷, with *traK*
500 found to be associated with industrialization status in five species (BOVA, BUNI, PDIS, PVUL,
501 PMER). Additionally, two gene families were annotated as relaxases, enzymes involved in the
502 mobilization of plasmids and other mobile elements. Other annotated genes included a putative
503 virulence factor and a lysozyme, both of which may play roles in microbial competition or host
504 interactions (Fig. 5B).

505

506 Together, these findings highlight the key role of horizontal gene transfer and mobile genetic
507 elements in shaping the accessory genome in response to host lifestyle^{48,49}. The recurrence of

508 industrialization-associated gene families across multiple species suggests that adaptation to
 509 industrialized and non-industrialized environments may be mediated, in part, by the acquisition
 510 and spread of genes involved in cell surface structure, conjugation, resistance, and species-
 511 species interactions.
 512



513
 514 **Figure 5 – Genes differentially enriched between host industrialized and non-industrialized**
 515 **lifestyles**

516 A. Gene enrichment analysis based on categorical and continuous levels of industrialization.
 517 Gene profiles were coded as presence/absence data and were correlated to host
 518 industrialization status encoded as a binary variable. Differential enrichment was tested
 519 while controlling for phylogeny. Significant hits (q -value < 0.05) are colored coded based
 520 on host lifestyle (purple: industrialized; green: non-industrialized). Non-significant genes
 521 are colored in grey. Top row: volcano plots showing all genes. The number of statistically
 522 significant genes is shown next to each species acronym. Significantly differentially
 523 enriched genes validated by measuring correlations with PC1 Lifestyle rather than
 524 industrialization status as a binary variable are shown in plain circles. Significant genes
 525 not validated with PC1 Lifestyle are shown as empty circles. Bottom row: top 10 most
 526 differentially enriched genes, for each lifestyle category. Genes with similar
 527 presence/absence profiles are collapsed into a single gene cluster. Gene labels indicate
 528 the number (n) of 90% gene families collapsed together.
 529 Data for 5 species are shown. These species harbor most of the significant hits. Data for
 530 the other 5 species is shown in Supp. Fig. 5.

531 B. Gene families (90% and 50% similarity gene clusters on top and bottom panels,
 532 respectively) with signals of differential enrichment across multiple species. Most gene
 533 families show convergent signals of differential enrichment based on host lifestyle
 534 (enrichment in one of two lifestyle categories across multiple species).
 535
 536

537 **Lifestyle-specific signals of selection at the genome and gene levels**

538

539 We next quantified gene-level signatures of selection to identify protein-coding genes under
540 positive selection in industrialized or non-industrialized host environments. To do so, we focused
541 on unicopy core protein-coding genes and calculated Ka/Ks ratios separately for strains from each
542 lifestyle category. On average, 1,724 genes were included in the Ka/Ks analysis per species
543 (range: 1,146 for *Bacteroides xylanivorans* to 2,065 for *Bacteroides ovatus*), with 87% of these
544 gene families having sufficient synonymous divergence (Ks) to be analyzed with confidence (see
545 Methods, Supp. Table 6).

546

547 Across species, 1.4% of genes showed evidence of positive selection ($Ka/Ks > 1$) in at least one
548 lifestyle category (Fig. 6A). The prevalence of positively selected genes varied substantially
549 between species, being highest in *P. dorei* (PDOR, 6.05%), *P. merdae* (PMER, 1.89%), and *P.*
550 *vulgatus* (PVUL, 1.81%), and lowest in *B. ovatus* (BOVA, 0.22%), *B. xylanivorans* (BXYL, 0.31%),
551 and *P. distasonis* (PDIS, 0.36%) (Fig. 6C). These differences are consistent with species-specific
552 variation in overall selection pressure, as reflected in median Ka/Ks values across genes, with
553 PDOR displaying the highest (0.28) and BXYL the lowest (0.09) (Fig. 6B). We also identified
554 numerous genes that were under positive selection specifically in either industrialized or non-
555 industrialized hosts, suggesting context-specific adaptive pressures (Fig. 6A & C-D).

556

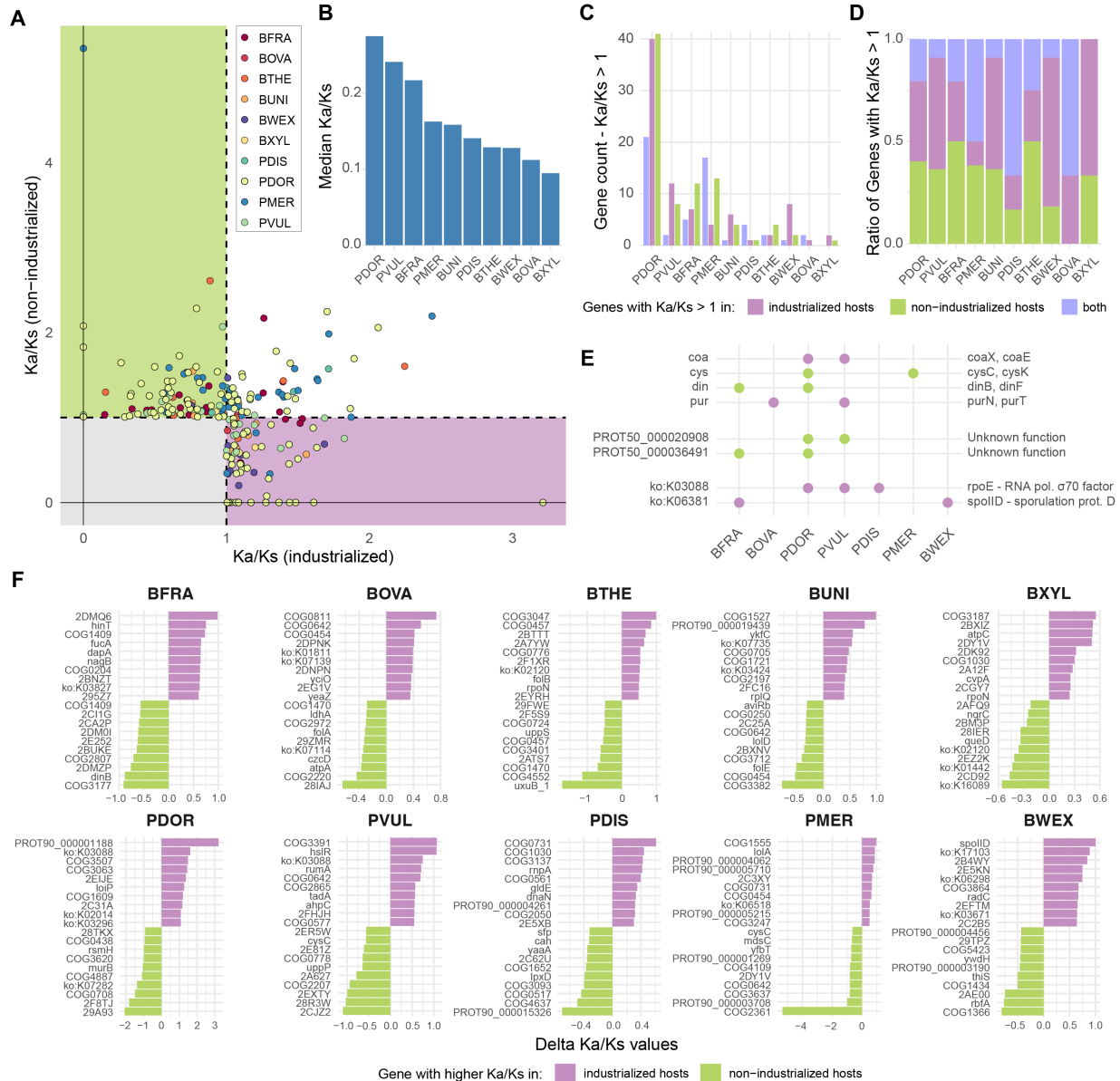
557 To collect additional evidence of positive selection associated with industrialization status, we
558 searched for convergent signals across species – specifically, gene functions under
559 industrialization-specific positive selection in multiple taxa. This analysis revealed two KEGG KOs
560 with repeated signatures of positive selection in strains from industrialized hosts (Fig. 6E). The
561 first, K03088, encodes rpoE, an extracytoplasmic function sigma factor from the sigma70 family,
562 which regulates the expression of stress response genes involved in maintaining cell envelope
563 integrity and responding to oxidative stress⁵⁰. K03088 genes are under positive selection in
564 strains of PDIS, PDOR, and PVUL in industrialized hosts, suggesting a lifestyle-associated
565 diversification of envelope stress responses in industrialized environments. The second, K06381,
566 groups genes with peptidoglycan lytic transglycosylase activity, including spoIID that is involved
567 in sporulation (Fig. 6E)⁵¹. Two genes within this orthologous group showed signals of positive
568 selection in industrialized hosts: spoIID itself in BWEX, and a SpoIID/LytB domain-containing
569 protein in BFRA. These findings suggest that peptidoglycan hydrolase functions – which play
570 roles in cell wall remodeling during sporulation in BWEX – may be adaptive targets in
571 industrialized gut environments.

572

573 We also identified six operons containing multiple genes under positive selection across species:
574 genes of the CoA biosynthesis (coa) (coaE in PVUL and coaX in PDOR) and of the purine
575 metabolism (pur) (purT in BOVA and purN in PVUL) operons are under positive selection in
576 industrialized hosts. Genes of the DNA damage response (din) (dinB in BFRA and dinF in PDOR)
577 and of the cysteine metabolism (cys) (cysK in PDOR and cysC in PMER) operons are under
578 positive selection in non-industrialized hosts (Fig. 6E).

579

580



581
582
583
584
585
586
587
588
589
590
591
592
593
594

Figure 6 – Lifestyle-specific signals of positive selection at the gene level

- A. Gene-level Ka/Ks values across species and host lifestyles (90% similarity gene families). For each lifestyle category and each gene, Ka/Ks values were computed for all pairs of codon-aligned gene sequences. Median Ka/Ks values are reported.
- B. Distribution of species-level median Ka/Ks values, aggregated across all genes.
- C. Counts of genes with median Ka/Ks values ≥ 1 (positive selection) across host lifestyle categories.
- D. Percentage of genes with median Ka/Ks values ≥ 1 across host lifestyle categories.
- E. Operons, 50% similarity gene clusters and KEGG KOs with convergent signals of positive selection across 2+ species.
- F. Top genes with highest absolute differences in median Ka/Ks values between industrialized and non-industrialized lifestyle categories.

595

596

597 To further investigate host lifestyle-driven selection, we ranked genes based on the difference in
598 Ka/Ks between strains from industrialized and non-industrialized hosts, regardless of whether
599 Ka/Ks exceeded 1 (Fig. 6F). This analysis revealed a set of genes with consistently elevated
600 selection pressure in one lifestyle category across multiple species (Supp. Fig. 6), many of which
601 are associated with stress response, metabolism, or virulence (Supp. Fig. 6). For instance, both
602 pepP and tpiA exhibited elevated Ka/Ks in industrialized hosts across eight species. pepP
603 encodes Xaa-Pro aminopeptidase P, a cytosolic metallo-exopeptidase involved in peptide
604 degradation, outer membrane vesicle (OMV) production, and bacterial virulence⁵². tpiA encodes
605 triosephosphate isomerase, a central enzyme in glycolysis essential for growth on glucose and
606 other glycolytic substrates⁵³. In several bacteria, tpiA has also been linked to virulence,
607 pathogenicity, and antibiotic resistance⁵⁴. Another gene, truB, showed elevated Ka/Ks in
608 industrialized hosts across seven species. truB encodes a tRNA pseudouridine synthase, which
609 contributes to ribosome stability and stress adaptation. truB has also been implicated in bacterial
610 virulence in pathogenic contexts⁵⁵.

611

612 Conversely, uppS and serC exhibited elevated Ka/Ks in non-industrialized hosts across eight and
613 seven species, respectively (Supp. Fig. 6). uppS encodes an undecaprenyl diphosphate
614 synthase, an essential enzyme in the biosynthesis of lipid carriers required for peptidoglycan
615 synthesis, cell wall assembly, and antibiotic resistance⁵⁶. serC encodes a phosphoserine
616 aminotransferase, which is critical for serine biosynthesis and supports metabolic flexibility,
617 particularly under nutrient-limited conditions^{57,58}.

618

619 Together, these results demonstrate that host industrialization status can exert selective
620 pressures on specific genes and pathways across multiple gut bacterial species.

621 **Parallel SNV-level adaptation to host industrialization across gut bacteria**

622

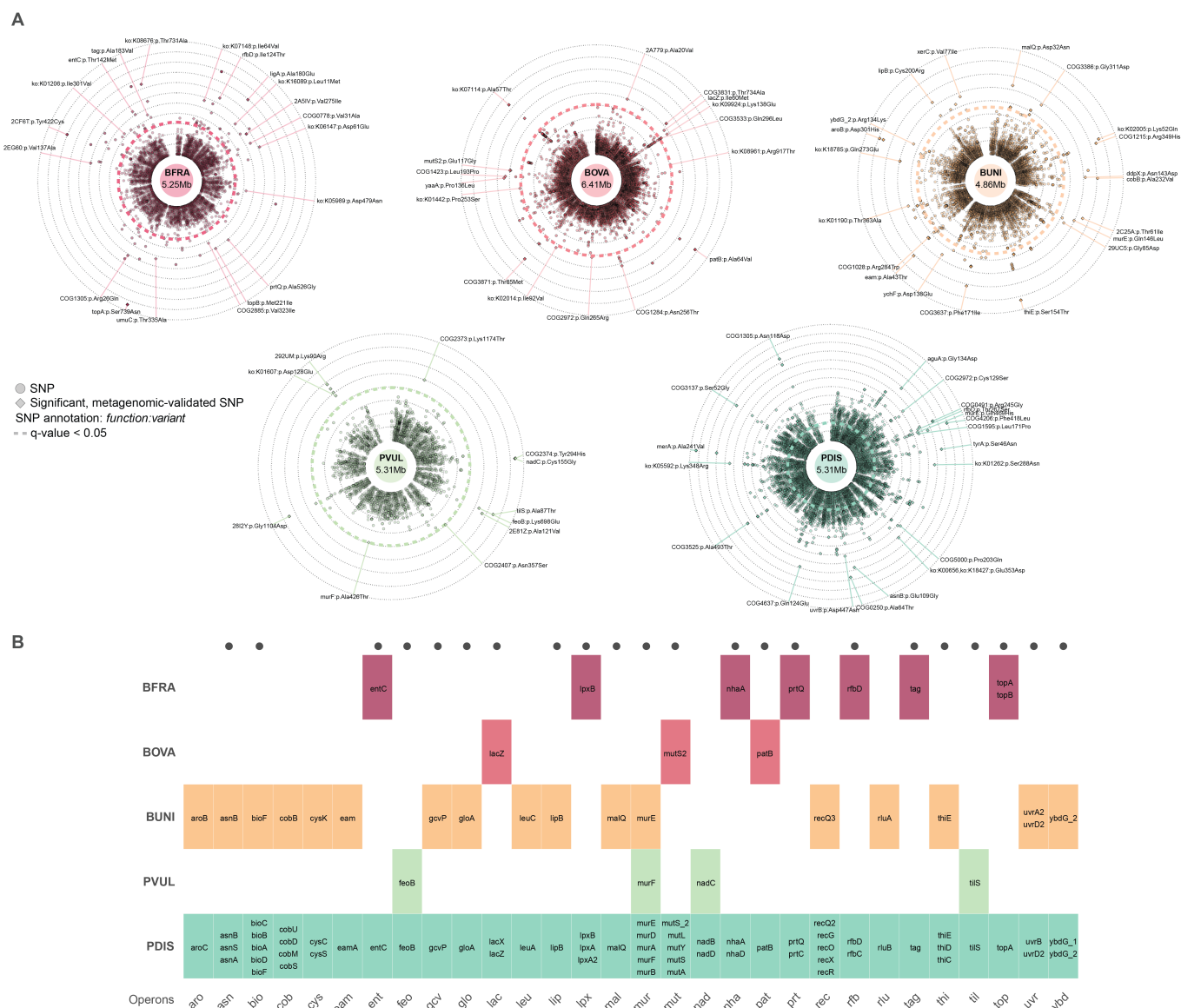
623 Next, we identified single nucleotide variants (SNVs) that are associated with host industrialization
624 status to go beyond the detection of positive selection at the level of genes, as individual variants
625 may also evolve in response to host environments. To investigate this, we analyzed non-
626 synonymous SNVs in the core genome, while accounting for phylogenetic structure (see
627 Methods). To limit the confounding effect of recombination, we excluded SNVs predicted to be
628 part of recombination regions⁵⁹ (see Methods). This analysis was restricted to the five species
629 with the largest number of genomes and StGBs (BFRA, BOVA, BUNI, PDIS and PVUL) to ensure
630 sufficient statistical power for detecting lifestyle-associated SNVs. To validate associations drawn
631 from the analysis of isolate genomes, we performed a secondary association analysis using SNVs
632 called from MAGs reconstructed from gut shotgun metagenomes of the entire GMbC cohort ($n =$
633 1,015 individuals, 12 countries and 35 localities spanning a range of a range of industrialization
634 levels and subsistence strategies), which we recently described⁶ (see Methods). We also
635 considered an additional layer of validation, using PC1 Lifestyle as a continuous response
636 variable for host industrialization. We then defined an SNV hit as “high-confidence” if (i) it showed
637 a consistent direction of effect with both isolate and MAG data, (ii) it achieved FDR-adjusted
638 significance ($p\text{-val} < 0.05$) in the isolate genome analysis, (iii) it achieved significance ($p\text{-val} <$
639 0.05) in the validation analysis with MAGs, and (iv) it achieved significance ($p\text{-val} < 0.05$) in the
640 validation analysis with PC1 Lifestyle (Supp. Table 7).

641

642 Across the five species tested, we found numerous high-confidence non-synonymous SNVs in
643 core genes that are significantly associated with host industrialization (Fig. 7A). The identification
644 of such SNVs suggests that convergent adaptation to host industrialization across distantly-
645 related strains is occurring within each of these species. Interestingly, 93%, 71%, 79%, 85% and
646 100% of hits detected with industrialization categories were validated with PC1 Lifestyle for BFRA,
647 BOVA, BUNI, PDIS and PVUL, respectively. Beyond showing robustness, these validations also
648 suggest that changes in non-synonymous SNV frequencies occur along a spectrum of host
649 lifestyle variation. Interestingly, genes containing these host lifestyle-associated SNVs are part of
650 operons involved in stress response (e.g. *mut*, *uvr*), nutrient acquisition and central metabolism
651 (e.g. *asn*, *nad*), host interaction and virulence (e.g. *ent*, *prt*), and transport and membrane
652 remodeling (e.g. *feo*, *mal*) (Fig. 7A & Supp Table 7).

653

654



655

656

657 **Figure 7 – Patterns of single nucleotide polymorphisms reveal genes and operons that**
 658 **undergo cross-species parallel evolution associated with host lifestyle**

659 A. Associations between single nucleotide variants (SNVs) and host lifestyle categories were
 660 calculated while accounting for phylogeny (see Methods). Associations were calculated
 661 for the 5 species with enough genome sample size to yield sufficient statistical power. Hits
 662 (q-value < 0.05) were cross-validated using GMbC shotgun metagenomes (n = 1,015, see
 663 Methods) that were sampled from diverse geographies worldwide, including those from
 664 which isolate genomes originate (reference). Hits validated by metagenomes are shown
 665 as diamonds and annotated.

666 B. Convergent signals of SNV-host lifestyle associations across bacterial species. Each tile
 667 represents an operon–species association and contains the names of genes within that
 668 operon that contain host-lifestyle associated SNVs. Operons are shown along the x-axis,

669 and species along the y-axis, cells are color-coded by species identity. Black points show
670 operons in which similar genes contain host-lifestyle associated SNVs across species. All
671 SNV data can be found in Supp. Table 7.

672
673

674 We then searched for genes and gene functions that harbor lifestyle-associated non-synonymous
675 SNVs in multiple species. As for our previous analyses on gene enrichment and Ka/Ks, the
676 rationale is that cross-species convergent evolution would provide even further compelling
677 evidence of adaptation to host environments shaped by industrialization beyond strain-level
678 convergent evolution. This analysis revealed 23 genes with lifestyle-associated non-synonymous
679 SNVs occurring in at least two different species (Fig. 7B). We further identified 30 operons
680 containing genes with lifestyle-associated non-synonymous SNVs across multiple species.
681 Remarkably, we found that the cob operon, which encodes enzymes for cobalamin (vitamin B12)
682 biosynthesis, harbors lifestyle-associated SNVs in PDIS and PUNI. This mirrors recent findings
683 from our analysis of GMbC gut metagenomes⁶, which identified cobalamin biosynthesis pathways
684 as significantly enriched in populations having increased levels of industrialization, after adjusting
685 for geographic, dietary, and host genetic confounders. These two independent lines of evidence,
686 found at the genomic and metagenomic levels, strongly suggest that vitamin B12 biosynthesis is
687 a key adaptive trait in response to industrialized host environments. Second, we observed
688 repeated signals in the mur operon, which is essential for murein (peptidoglycan) biosynthesis, a
689 core component of the bacterial cell wall. Specifically, lifestyle-associated non-synonymous SNVs
690 were identified in murE in BUNI, murF in PVUL, and murA, murB, murD, murE, and murF in PDIS.
691 This points to a potential role of cell wall remodeling and integrity in adaptation to differing host
692 lifestyles. Third, we found that two operons involved in lipopolysaccharide (LPS) biosynthesis (lpx,
693 encoding genes for lipid A biosynthesis and rfb, encoding genes for O-antigen biosynthesis) also
694 harbor host lifestyle-associated SNVs in multiple species. Given the role of LPS in host immune
695 recognition and membrane structure, these findings suggest that modulation of LPS structure may
696 be an important contributor to adaptation to host lifestyle-specific environments.

697

698 Together, these findings provide evidence for SNV-level parallel evolution across gut bacteria,
699 highlighting shared functional targets that may mediate adaptation to host lifestyle and
700 industrialization.

701
702
703
704
705

706 Discussion

707 In this study, we provide evidence that human lifestyle, particularly factors associated with
708 industrialization status, shape the genomic evolution of bacterial commensals. We show that
709 strains from industrialized hosts harbor larger proteomes and exhibit increased pangenome
710 fluidity. We found that recently acquired genes during bacterial evolution, explained by HGT, are
711 driving these phenomena. These findings confirm that HGT is a primary mechanism for rapid
712 evolutionary innovation in the human gut microbiome, and are in line with our previous discovery
713 that frequencies of HGT are elevated among microbes co-occurring in industrialized hosts ²¹.
714 Interestingly, a recent study also showed that HGT facilitates adaptative selective sweeps of
715 *mdxEF* genes, involved in maltodextrin metabolism, among industrialized populations ⁴⁸. We
716 further accumulated multiple lines of evidence of convergent evolution of genomic features with
717 respect to host industrialization status and lifestyle (PC1 Lifestyle), including the differential
718 enrichment of genes, gene-level positive selection, and non-synonymous SNV frequency. Such
719 convergent patterns could be observed at different taxonomic resolutions, both between strains
720 and across species. These parallel genomic signatures are consistent with adaptation to the
721 ecological conditions associated with industrialized or non-industrialized environments, including
722 exposure to different commensal and environmental microbes, dietary profiles and medication
723 usage ⁶⁰. The functional categories involved in these changes include traits that are relevant for
724 ecological adaptation. In particular, genes under selection or containing lifestyle-associated SNVs
725 are involved in stress response (e.g., *rpoE*, *uvr*, *mut*), cell envelope remodeling (e.g., *mur*, *lpx*,
726 *rfb*), central metabolism (e.g., *cob*, *nad*, *asn*), and host interaction or virulence-related functions
727 (e.g., *prt*, *ent*) (Fig. 5 & 6).

728
729 Lastly, several of the adaptive features that we identified may have important implications for host
730 physiology and health. Notably, genes involved in vitamin B12 biosynthesis harbor lifestyle-
731 associated non-synonymous SNVs across multiple bacterial species. Several vitamin B12
732 biosynthesis pathways were also found to be differentially abundant in the metagenome of the
733 broader GMbC cohort (n = 1,015) ⁶, with higher abundance among populations more exposed to
734 industrialized environments. This suggests that cobalamin biosynthesis constitutes a key adaptive
735 trait in industrialized gut ecosystems. Similarly, we observed multiple signals of adaptation in
736 lipopolysaccharide (LPS) biosynthesis pathways – including *wecA*, *lpx*, and *rfb* genes (Figs. 5 &
737 7) – which are critical for bacterial membrane structure, bacteria-bacteria interactions and host
738 immune recognition. These findings indicate that LPS remodeling may not only contribute to
739 microbial adaptation in industrialized hosts, but could also influence host immune and
740 inflammatory responses ⁶¹.

741 Limitations

742 Our analysis focused on ten bacterial species. While these species are prevalent and functionally
743 significant components of the human gut microbiome, further investigations that broadens
744 taxonomic and host representation will be necessary to assess the generalizability of our findings.
745 In addition, even though we discovered several associations between host industrialization status
746 and microbial genomic features, it is still unclear which environmental factors, such as dietary
747 composition, medication use, sanitation, or other exposures related to lifestyle, are responsible
748 for these correlations. Future research combining functional validation, experimental evolution,
749

750 and fine-scale environmental variable measurement will be essential to identify and separate the
751 selective pressures that influence the evolution of microbial genomes in response to variations in
752 the industrialization of host environments.

753

754 Finally, the evolutionary mechanisms driving the expansion of proteomes in strains from
755 industrialized hosts remain unclear. Although we find signatures of adaptation shaping the
756 distribution of individual genes, the overall proteome expansion may result from the relaxation of
757 purifying selection, allowing the accumulation of slightly deleterious genes and mobile elements.
758 Future studies involving population structure reconstruction, effective population size estimation,
759 the study of mutational and recombination processes influencing allele frequencies, and fitness
760 experiments will be required to distinguish the relative contributions of adaptive versus non-
761 adaptive processes in influencing the evolution of the bacterial genome across lifestyles.

762

763

764

765

766

767

768

769 **Methods**

770

771 **Participant recruitment and collection of biospecimens**

772 To culture and isolate the additional 1,841 GMbC bacterial strains generated in this study, we
773 selected 20 participants from the GMbC cohort (cross-sectional population cohort of healthy
774 adults that we recently described⁶). These individuals were chosen to represent a broad range
775 of lifestyle, urbanization levels, and geographic origins (see below), and none had taken
776 antibiotics or antiparasitic treatments recently. All participants were asymptomatic for infectious
777 or chronic diseases at the time of enrollment. The 20 participants were from Cameroon, Ghana,
778 Malaysia, Nigeria, and Rwanda. This new set of isolate genomes builds upon our initial GMbC
779 collection of 4,140 isolates²¹, bringing the total to 6,000 isolates from 56 participants spanning
780 nine countries, including the USA, Canada, Finland, and Tanzania. Written informed consent was
781 obtained from all participants, using translations in local language when appropriate. Research &
782 ethics approvals were obtained from the MIT IRB (protocol #1612797956) and from the Ethics
783 commission of the Medical Faculty of Kiel University (Studie D 511/24). Permits were also
784 obtained in each sampled country prior to the start of sample collection, from the following ethics
785 committees:

- 786 • Cameroon: Comité National d’Ethique de la Recherche pour la Santé Humaine, protocol
787 #2017/05/901/CE/CNERSH/SP;
- 788 • Ghana: Cape Coast Teaching Hospital Ethical Review Committee, protocol
789 #CCTHERC/RS/EC/2016/3 and Committee on Human Research, Publication and Ethics
790 of the Komfo Anokye Teaching Hospital, protocol #CHRPE/AP/398/18;
- 791 • Malaysia: Universiti Malaya Medical Research Ethics Committee, MREC ID No.: 2018219-
792 6033;
- 793 • Nigeria: National Health Research Ethics Committee of Nigeria, protocol
794 #NHREC/01/01/2007-29/04/2018.
- 795 • Rwanda: National Ethics Committee, protocol IRB 00001497 of IORG0001100

796

797 Participants self-collected a fresh fecal sample using sterile containers, which were returned to
798 GMbC scientists for on-site processing within 3 hours after defecation. Raw stool was diluted 1:5
799 in a 25% pre-reduced anaerobic glycerol solution containing acid-washed glass beads for
800 homogenization, then aliquoted into 2 mL cryogenic tubes. Aliquots in cryoprotectant were flash-
801 frozen either in liquid nitrogen (−196 °C) using a cryoshipper tank or stored at −80 °C. An
802 additional 1–2 g of stool was preserved in RNAlater for DNA stabilization and sequencing. All
803 samples were subsequently shipped to MIT for further processing and storage.

804

805 **Host metadata and industrialization status**

806 Lifestyle metadata were collected for all GMbC participants and analyzed using dimensionality
807 reduction techniques⁶. Briefly, we used the Human Development Index (HDI), as described
808 previously²¹, as a measure to determine the industrialization status of sampled populations.
809 Populations from localities with HDI values below the national median (0.739 in 2022) were
810 classified as non-industrialized, while those with values above this threshold were classified as
811 industrialized. We also collected a range of lifestyle variables including, e.g., population density,
812 main subsistence strategies, type of floor in households, access to electricity, household size,

813 contact with animals, frequency of physical activities per week and source of drinking water. A full
814 list of lifestyle variables, along with dietary and medication data for the participants used for
815 bacterial culturing is provided in Supp. Table 1.

816 We used the PCAmix function from the PCAmixdata R package to perform a principal component
817 analysis of lifestyle variables (Fig. 1B). PC1 Lifestyle strongly correlates with industrialization
818 status, along with other key factors such as access to electricity, main subsistence strategy, floor
819 type and type of drinking water (Fig. 1B). PC1 Lifestyle is used as a continuous proxy for
820 industrialization.

821

822 **Culturing and isolation of gut bacteria**

823 We used the same culturing procedures as used previously to build the BIO-ML¹² and the GMbC
824²¹ isolate collections. Fecal samples were processed anaerobically (5% Hydrogen, 20% Carbon
825 dioxide, balanced with Nitrogen) and diluted in pre-reduced PBS (with 0.1% L-cysteine
826 hydrochloride hydrate). Samples were plated onto pre-reduced agar plates and incubated
827 anaerobically at 37C for 7 to 14 days. Both general (nonselective) and selective media were used
828 to culture diverse groups of organisms (Supp Table X). After incubation, bacteria were isolated
829 and purified through two rounds of picking and re-streaking at 37C. One colony was then
830 inoculated in liquid media. After 2 days of anaerobic incubation at 37C, the taxonomy of the isolate
831 was identified using 16S rRNA gene Sanger sequencing (starting at the V4 region). We first
832 amplified the full 16S rRNA gene by PCR (27f 50-AGAGTTTGATCMTGGCTCAG-30 - 1492r 50-
833 GGTTACCTTGTTACGACTT-30) and then generated a 1kb long sequence by Sanger reaction
834 (u515 50-GTGCCAGCMGCCGCGGTAA-30). Isolates were then stored at -80C in a pre-reduced
835 cryoprotectant glycerol buffer. See Supp Table X for the list of culturing media, selection
836 treatments and isolate metadata.

837

838 **DNA extraction, library prep and sequencing of isolate genomes**

839 We used the same procedures used to generate the first set of GMbC isolate genomes²¹. Briefly,
840 whole genome DNA from isolates was extracted with the DNeasy UltraClean96 MicrobioaKit
841 (QIAGEN), following manufacturers' protocols. Genomic DNA libraries were built from 1.2ng of
842 extracted DNA using the Nextera DNA Library Preparation kit (Illumina), following the
843 manufacturer's protocol, with reaction volumes scaled accordingly. Isolated libraries were pooled,
844 with insert size and concentration of each pooled library being determined using an Agilent
845 Bioanalyzer DNA 1000 kit (Agilent Technologies). Paired-end (2x150bp) reads sequencing was
846 performed using an Illumina NextSeq 500 instrument (Illumina Inc) at the Broad Institute.

847

848 **Assembly and quality estimation of GMbC isolate genomes and MAGs**

849 Isolate genome assemblies were reconstructed as follows. We first used cutadapt v1.12⁶² (with
850 parameters -a CTGTCTCTTAT -A CTGTCTCTTAT) and Trimmomatic v0.36⁶³ (with parameters
851 PE -phred33 LEADING:3 TRAILING:3 SLIDINGWINDOW:5:20 MINLEN:50) on demultiplexed
852 paired-end short reads to remove barcodes and Illumina adapters and for base quality filtering.
853 We then used SPAdes v.3.9.1⁶⁴ (with parameter-careful) for de novo assembly of reads into
854 contigs. To iteratively improve genome assemblies, we used SSPACE v3.0⁶⁵ and GapFiller v1-
855 10⁶⁶ to scaffold contigs and to fill sequence gaps (with default parameters). We removed scaffolds
856 smaller than 1kb. All reads were aligned back to the assembly to compute genome coverage

857 using BBmap v37.68 (<https://jgi.doe.gov/data-and-tools/bbtools/>) and the covstats option (with
858 default parameters).

859 MAGs were reconstructed from $n = 1,015$ shotgun metagenomes of the GMbC cohort, as
860 described in our recent work ⁶, resulting in a total of $n = 24,163$ quality-filtered MAGs (see below).
861 Among these, 1,424 MAGs originated from the donors used for generating GMbC isolate
862 genomes (from the same original stool sample). Briefly, quality filtered short read shotgun
863 metagenomic data was assembled into contigs using Megahit ⁶⁷. Contigs were binned using a
864 multi-binning and refinement strategy in MAGScoT ⁴⁰, using MaxBin2 ⁶⁸, Metabat2 ⁶⁹, CONCOCT
865 ⁷⁰ and VAMB ⁷¹ binning tools.

866 Genome completeness and contamination were estimated using CheckM and MAGScoT ^{40,72}.
867 Standard thresholds for completeness ($>50\%$) and contamination ($<10\%$) were used to filter out
868 low-quality genomes. Genome quality was defined as $Q = \text{Completeness} - 0.5 * \text{Contamination}$.
869 We considered genomes with $Q > 0.7$ as high-quality genomes, and genomes with $0.7 > Q > 0.5$
870 as medium-quality genomes. All GMbC isolate genomes are high-quality genomes, with median
871 $Q = 0.98$, min $Q = 0.88$, max $Q = 1$. Assembly and genome quality statistics for GMbC isolate
872 genomes and MAGs can be found in Supp. Table 1.

873

874 **External collections of human gut bacterial isolate genomes and MAGs**

875 We included 3,632 human gut bacterial isolate genomes from the BIO-ML collection, which we
876 generated in a previous study ¹². These isolates were obtained from participants with
877 industrialized lifestyles recruited in the Boston (MA, USA) area and were processed using the
878 same culturing, sequencing, and assembly protocols described above for the GMbC isolate
879 genomes. MAGs were also reconstructed from the 11 BIO-ML participants used for culturing and
880 isolating bacteria, using the same pipeline as for GMbC MAGs. For comparative analyses, we
881 also included 4,644 representative MAGs from the UHGGv2 collection of human gut MAGs ^{24,36}.

882

883 **GMbC shotgun metagenomes**

884 We generated shotgun gut metagenomic data for GMbC cross-sectional population cohort ($n =$
885 $1,015$ healthy adult participants). This includes the participants whose sample we used in this
886 study to culture and isolate gut bacteria. We present the shotgun metagenomic data in our
887 recently published work ⁶. Briefly, the data were sampled from 12 countries and 35 localities,
888 across a range of industrialization levels and subsistence strategies. Briefly, DNA was extracted
889 from stool samples stored in RNAlater using the MoBio Powersoil 96 kit, and sequencing libraries
890 were prepared with the Nextera XT kit. Paired-end 150 bp sequencing was performed on the
891 Illumina NovaSeq S4 platform, generating a median of 21.7 million reads per sample. GMbC and
892 BIO-ML MAGs were reconstructed as described above (section "Assembly and quality estimation
893 of GMbC isolate genomes and MAGs"). GMbC and BIO-ML MAGs and isolate genomes were all
894 clustered into reference species-level genome bins (SGBs) using a multi-step workflow based on
895 ANI clustering with dRep and fastANI, resulting in a final non-redundant set of 2,379 SGBs, 480
896 of which contained isolate genomes ⁶. SGB representatives were taxonomically annotated with
897 GTDB-Tk (v2.3) and the GTDB reference database (rel214). They were used as references for
898 abundance estimation with Salmon in metagenome mode ⁷³. Relative abundances were
899 expressed in TPMs, filtered to remove low-abundance SGBs ($<1,000$ reads and <250 TPMs), and
900 transformed using CLR with a pseudocount of 1. Fig. 1B shows the microbiome alpha and beta

901 diversity of GMbC and BIO-ML participants from whom bacterial isolates were obtained. Alpha
902 diversity was calculated using Faith's PD index, while beta diversity was assessed using the
903 unweighted UniFrac dissimilarity implemented in the GUniFrac R package ([https://cran.r-](https://cran.r-project.org/web/packages/GUniFrac/index.html)
904 [project.org/web/packages/GUniFrac/index.html](https://cran.r-project.org/web/packages/GUniFrac/index.html)). Both metrics were computed using a
905 phylogenetic tree of the reference SGBs based on the GTDB single-copy marker gene alignments
906 and the GTDBtk "infer" workflow.

907

908 **SGB and StGB reconstruction from isolate genomes and MAGs**

909 We clustered genomes into species-level genome bins (SGBs) by combining GMbC and BIO-ML
910 isolate genomes with GMbC MAGs. To ensure high-quality clustering, we implemented an
911 iterative workflow that optimizes cluster resolution and quality. First, GMbC MAGs were clustered
912 within each geographic location at 97% average nucleotide identity (ANI) using dRep (v3.4.0).
913 For each cluster, we selected the MAG with the highest quality score as its representative. These
914 representative MAGs were then classified as either high-quality (HQ; score > 0.7) or medium-
915 quality (MQ; $0.5 \leq \text{score} \leq 0.7$). Next, HQ MAG representatives were combined with all GMbC
916 and BIO-ML isolate genomes and clustered into SGBs (95% ANI). Each MQ MAG representative
917 was then compared to the HQ SGB representatives using fastANI (v1.33) and assigned to an
918 SGB if it shared at least 95% ANI. MQ MAGs without matches at the required threshold were
919 reclustered using dRep to form MQ SGBs, and we conserved bins only if they contained more
920 than one genome. Finally, HQ and MQ clusters were merged to generate the final set of non-
921 redundant SGBs. GMbC and BIO-ML isolate genomes were also clustered at 99% ANI to form
922 strain-level genome bins (StGBs). For the evolutionary genomic analyses of the 10 species
923 presented in Fig. 3-7, StGBs containing genomes from multiple donors were further split by donor
924 to ensure that each StGB represented a donor-specific strain.

925

926 **Comparison of GMbC isolates with UHGG MAGs, BIO-ML isolates and GMbC MAGs**

927 SGB and StGB representative genomes were compared to the UHGGv2 representative MAGs (n
928 = 4,644) using fastANI. A genome was classified as "included in UHGG" if it shared >95% ANI
929 with a UHGG representative genome. GMbC isolate genomes were considered to be represented
930 in GMbC MAGs or BIO-ML isolates if they clustered within the same SGB or StGB.

931

932 **Identification of Isolate-MAG pairs of genomes**

933 Isolate-MAG pairs from the same donor and same species-level taxonomy were identified by
934 matching MAGs and isolate genomes from the same individual that clustered within the same
935 SGB. If multiple isolate genomes from a given donor were present within an SGB, we calculated
936 pairwise ANI values between each isolate and the corresponding MAG. The isolate genome with
937 the highest ANI to the MAG was selected as the representative genome for downstream pairwise
938 comparisons of genomic features (Fig. 2).

939

940 **Gap-seq metabolic features**

941 Genome-scale metabolic models were reconstructed using gapseq (v1.2, commit 2dfa8c80)⁷⁴ for
942 all paired MAGs and isolate genomes (Fig. 2). The reconstruction workflow using gapseq
943 consisted of 5 steps. First, all bacterial pathways from the MetaCyc and gapseq database were
944 retrieved with 'gapseq find' and options '-p all -t Bacteria'. Second, all cross-membrane metabolite

945 transporters were recovered using ‘gapseq find-transport’. Third, draft metabolic networks were
946 reconstructed from the two previous steps using ‘gapseq draft’. Fourth, an anoxic growth medium
947 for the given organism is predicted from the draft metabolic network using ‘gapseq medium’ and
948 the option ‘-c "cpd00007:0"’, following our previous study ⁷⁵. Finally, gaps in the draft metabolic
949 network were filled in order to enable growth (i.e. formation of biomass), assuming the growth
950 medium predicted at the previous step. This step was performed using the ‘gapseq fill’ module.
951 For all models, the number of reactions, metabolites, and genes that are linked to reactions or
952 transporters are counted for the draft network reconstruction and for the gap-filled network.

953

954 **Functional profiling of ARGs, VFs, CAZymes, MGEs and MGE machineries in isolate-MAG** 955 **pairs**

956 Antibiotic resistance genes and virulence factors were identified in genome assemblies using
957 Abricate (<https://github.com/tseemann/abricate>), with annotations based on the NCBI
958 AMRFinderPlus database ⁷⁶ for resistance genes and the VFDB database ⁷⁷ for virulence factors.
959 CAZyme genes were detected using dbCAN3 (https://github.com/linnabrown/run_dbcan) ⁷⁸. We
960 used MacSyFinder (version 2.0) ⁷⁹ to profile MGEs and MGE machineries of three key systems:
961 type secretion systems (TXSS), type IV filaments (TFF), and conjugative transfer systems (Conj).
962 We used ISEScan (version 1.7.2.3) ⁸⁰ to identify insertion sequences (IS). We used geNomad
963 (v1.7.4) ⁸¹ to detect prophages. Following a recently established protocol ⁸², we first used checkV
964 (v1.5) ⁸³ to measure phage sequence quality, retaining only prophage sequences with >50%
965 completeness (including medium-quality, high-quality and complete sequences) for downstream
966 analysis. Prophage sequences were then clustered based on average nucleotide identity, with
967 the longest prophage contig in each cluster selected as the representative sequence.

968

969 **Inference of horizontal gene transfers from isolate genomes and MAGs**

970 We detected HGTs on the 147 pairs of MAG and isolate genomes, identifying HGTs among MAGs
971 and isolate genomes separately. We focused on HGTs occurring between bacterial species. The
972 147 genomes cluster into 87 different SGBs, which constitutes 3,741 species pairs. To detect
973 HGTs, we used methods that we implemented in previous studies ^{21,41} that rely on the Blast-based
974 detection of blocks of DNA that are shared by two genomes of different species. We retained blast
975 hits with 100% similarity and that are larger than 500bp, to focus on the most recent HGTs ²¹. We
976 analyzed HGT sequences that have a relative read coverage higher than 0.2 compared to the
977 average genome coverage in at least one of the two compared genomes ²¹. We considered a
978 HGT having occurred between two species when we could identify at least one genome pair with
979 at least one blast hit.

980

981 **Species selection**

982 We ranked species in the combined GMbC and BIO-ML isolate genome collection based on the
983 total number of distinct StGBs sampled, as well as their representation across industrialized and
984 non-industrialized hosts. For downstream genomic analyses, we selected the top 10 species
985 meeting these criteria (Fig. 1G), each represented by a minimum of eight StGBs per host
986 industrialization category (Supp. Table X). Included species are *Bacteroides fragilis* (BFRA),
987 *Bacteroides ovatus* (BOVA), *Bacteroides thetaiotaomicron* (BTHE), *Bacteroides uniformis*
988 (*BUNI*), *Bacteroides xylanisolvens* (BXYL), *Phocaeicola dorei* (PDOR), *Phocaeicola vulgatus*

989 (PVUL), *Parabacteroides distasonis* (PDIS), *Parabacteroides merdae* (PMER), and *Blautia A*
990 *wexlerae* (BWEX). We purposefully excluded *Escherichia coli*, as our focus was on abundant
991 members of the gut microbiome that have not yet been extensively characterized – unlike *E. coli*,
992 whose global genomic diversity has been well studied ⁸⁴.

993

994 **Variant calling and annotation**

995 For the ten species, we considered all GMbC and BIO-ML isolate genomes with >95%
996 completeness, based on MAGScoT estimates. A reference genome per species was chosen
997 based on the best quality score Q (see above) and the highest L50 value. Contigs were joined
998 into a single fasta sequence with a stretch of 100 ambiguous bases between contigs to create an
999 artificial continuous genome for downstream read mapping and variant calling. Short read data of
1000 all isolate genomes were mapped against their respective species reference genome using snippy
1001 with default parameters, requiring a coverage ≥ 10 -fold for the variant calling. MAGs derived
1002 from metagenomic data of the broader GMbC cohort ⁶ that are of the same species were also
1003 included. We selected MAGs that have >97% average ANI (using skani) to any of the isolate
1004 genomes in each species, and that have low contamination estimates (<10%) data. MAGs were
1005 processed with snippy in the “contigs” mode (--ctgs flag) (<https://github.com/tseemann/snippy>).
1006 Isolate and MAG-derived snippy variant calls were combined in a single whole-genome alignment
1007 using the ‘snippy-core’ command. The alignment was then cleaned using the ‘snippy-
1008 clean_full_aln’ utility script provided by the snippy developer. To discriminate gaps within the
1009 alignment from gaps at the end of contigs, gaps of length ≤ 10 bp were masked. We discarded
1010 genomes that had more than 40% ambiguity calls or gaps from the final snippy output. The
1011 alignment was converted to a VCF file of all single-nucleotide variants (SNVs) using ‘snp-sites’
1012 utility script of snippy. Multi-allelic variants were decomposed to individual entries in the VCF file
1013 using vt decompose (<https://github.com/atks/vt>). Variant effects were predicted using snpEff ⁸⁵.

1014

1015 **Recombination-aware reconstruction of StGB phylogenies**

1016 A reference genome per StGB and per individual was selected using the same criteria as for the
1017 selection of species-level reference genomes (see previous section). Reference genomes were
1018 extracted from the whole genome alignment reconstructed with snippy and used for phylogeny
1019 reconstruction using IQ-TREE2 and a generally time reversible (GTR) substitution model ⁸⁶. We
1020 then used Gubbins ⁵⁹ to identify and mask regions of homologous recombination. The masked
1021 alignment was then reused in IQ-TREE2 (GTR model) for phylogenetic reconstruction.

1022

1023 **Gene calling and pangenome reconstruction**

1024 Protein-coding gene sequences from all isolate genomes of the ten selected species were
1025 predicted using PROKKA (default parameters) ⁸⁷. Coding sequences (CDSs) were then clustered
1026 into gene families at 90% sequence identity using MMseqs2 (v13.45111) and the ‘easy-lincludt’
1027 module with the options ‘--cov-mode 1 -c 0.8 --kmer-per-seq 80’ ⁸⁸. Representative sequences for
1028 each gene family were functionally annotated using EggNOG-mapper (v2.1.3; database version
1029 220425) ⁸⁹. Within each species, gene families were subsequently categorized into three groups
1030 based on their prevalence across genomes: (soft-)core genome (>95% of genomes), accessory
1031 genome (5–95%), and cloud genome (<5%).

1032

1033 **Quantifying mutational selective pressures**

1034 Species-level single-copy core gene families (SCGs) were defined as being present in 95% of all
1035 isolate genomes of a species, and found in each individual StGB. SCG protein sequences of the
1036 representative StGB genomes were aligned using MAFFT in '--auto' mode⁹⁰. Protein-level
1037 alignments were translated back to nucleotidic sequences to obtain codon-level alignments, which
1038 were subjected to pairwise Ka and Ks value calculations using the kaks() function of the 'seqinr'
1039 R package. Sequence pairs with at least one synonymous and one non-synonymous mutation
1040 were included, while others were assigned a -1 value and excluded from the analysis as
1041 recommended in PAML. For each gene family, mean Ka/Ks values were calculated for each host
1042 industrialization category. Mean Ka/Ks values for which > 90% of pairwise comparisons were
1043 excluded due to lack of mutations were marked as 'low confidence'.

1044

1045 **Calculation of the pangenome fluidity**

1046 We calculated pangenome fluidity as described previously^{91,92}. Briefly, pangenome fluidity
1047 measures the average proportion of genes that are not shared between pairs of genomes from
1048 the same species. For each species, we computed fluidity separately for strains isolated from
1049 industrialized and non-industrialized hosts, using the representative genome of each StGB.
1050 Specifically, we calculated the pairwise proportion of non-shared gene families between all StGB
1051 pairs within each lifestyle category. To compare the distribution of non-shared gene proportions
1052 between lifestyles, we performed Wilcoxon rank-sum tests.

1053

1054 **Gene tree – species tree phylogenetic reconciliations**

1055 We used the species tree-aware gene tree reconciliation method, ALE⁹³ implemented in AleRax
1056⁴⁴ (commit version: 8705a60, 2025-06-02) to harvest signals of gene evolutionary events (gene
1057 duplications, transfers and losses) from all the gene families of the 10 selected species. We used
1058 the recombination-aware strain-level StGB phylogenies (see above) as species trees, and all
1059 individual gene trees of the gene families containing at least 4 representative genomes. The
1060 species tree was midpoint-rooted using ete3⁹⁴. The nucleic acid sequences were ordered based
1061 on the 90% sequence identity clustering into per-family sequence files and deduplicated (keeping
1062 only the gene families which still had at least 4 unique sequences for downstream analysis). The
1063 number of gene families per species distributes as follows: BFRA: 4095; BOVA: 6238; BTHE:
1064 4436; BUNI: 4906; BWEX: 3596; BXYL: 5047; PDIS: 5042; PDOR: 2968; PMER: 3340; PVUL:
1065 4988. Gene families were aligned using MAFFT (with L-INS-i option)⁹⁰ with default settings. From
1066 the alignments, gene trees were inferred under the GTR+G+F model using IQ-TREE2⁸⁶ and
1067 generating 1000 ultrafast bootstrap samples⁹⁵ with the '-bnni' option. AleRax analysis was
1068 performed on these ultrafast bootstrap gene tree samples of all the gene families using the
1069 GLOBAL model parametrization (i.e. the DTL-rates are shared among families and jointly
1070 estimated), the UndatedDTL reconciliation model and taking 1000 reconciled gene trees along
1071 the species tree. Based on the reconciled gene trees, AleRax reports the number of evolutionary
1072 events along the branches and leaves of the species tree. Wagner-parsimony⁹⁶ was used on the
1073 species tree to deduce host industrialization status at internal nodes based on existing status at
1074 terminal nodes. Evolutionary events could then be associated with industrialization state on a per
1075 branch level (Fig. 4A). All reconciliation data can be found in Supp Table 4.

1076

1077 **Gene Family enrichment analysis**

1078 Accessory gene families (present in 5–95% of genomes) were included in the gene enrichment
1079 analysis if they were found in at least five StGBs and absent from at least five others. Gene
1080 distribution data were encoded as presence/absence data. To detect differentially enriched genes
1081 while accounting for phylogenetic structure, we used a phylogeny-aware logistic regression model
1082 implemented in the phylolm R package⁹⁷. Host industrialization status (industrialized vs. non-
1083 industrialized) was used as a binary explanatory variable, and we additionally tested associations
1084 using PC1 Lifestyle, a continuous proxy of industrialization derived from multidimensional lifestyle
1085 data (see section above). Significant associations identified with the binary industrialization
1086 variable were validated using PC1 Lifestyle and a phylogeny-aware linear regression model in
1087 phylolm. Recombination-aware phylogenies at the StGB level were used in all phylolm models.
1088 All genomes within each species were included in the regression analyses, with an arbitrary
1089 branch length equal to 1% of the median cophenetic distance added to non-representative isolate
1090 genomes to preserve tree topology.

1091

1092 **Detection of bacterial single nucleotide variants (SNVs) associated with host**
1093 **industrialization status**

1094 For SNV–host industrialization associations, we analyzed biallelic and non-synonymous variants
1095 located within single-copy core gene families (defined as present in >99% of genomes), retaining
1096 only variants with a minor allele frequency >5%. Variants located in predicted recombination
1097 regions were masked prior to analysis. Associations were tested using a logistic regression model
1098 that accounts for phylogenetic structure (‘phylolm’ R package). Host industrialization was
1099 considered both as a binary response variable (industrialized vs. non-industrialized) and using
1100 PC1 Lifestyle as a continuous proxy, consistent with the approach used in the gene enrichment
1101 analyses. The latter served as a validation for significant associations identified using the binary
1102 variable.

1103 To further validate SNV-level associations, we examined corresponding significant SNVs in the
1104 metagenomic data of the broader GMbC cohort by mapping metagenomes against the reference
1105 genome of each species (see section “Variant calling and annotation” above). Associations with
1106 metagenome-derived SNVs were also tested using logistic regression against host
1107 industrialization status as a binary variable.

1108

1109

1110

1111

1112

1113 **Data availability**

1114 Short read data and assemblies of GMbC isolate genomes generated in this study will be made
1115 available online on the dbGaP server (Study ID: 38715; Accession: phs002235.v1.p1; Accession:
1116 phs002205.v1.p1) upon publication of the article. GMbC metagenomes used in this study and
1117 published in our recent study ⁶ will be made available at the same dbGaP study.

1118

1119 **Code availability**

1120 Scripts and command lines used to process the data will be made available at
1121 <https://github.com/MMmicrobiome-Lab> upon publication of the article.

1122

1123

1124

1125 References

- 1126
- 1127 1. Schnorr, S.L., Candela, M., Rampelli, S., Centanni, M., Consolandi, C., Basaglia, G., Turrioni, S.,
- 1128 Biagi, E., Peano, C., Severgnini, M., et al. (2014). Gut microbiome of the Hadza hunter-gatherers.
- 1129 Nat. Commun. 5, 3654. <https://doi.org/10.1038/ncomms4654>.
- 1130 2. Carter, M.M., Olm, M.R., Merrill, B.D., Dahan, D., Tripathi, S., Spencer, S.P., Yu, F.B., Jain, S., Neff,
- 1131 N., Jha, A.R., et al. (2023). Ultra-deep sequencing of Hadza hunter-gatherers recovers vanishing gut
- 1132 microbes. Cell 186, 3111–3124.e13. <https://doi.org/10.1016/j.cell.2023.05.046>.
- 1133 3. Smits, S.A., Leach, J., Sonnenburg, E.D., Gonzalez, C.G., Lichtman, J.S., Reid, G., Knight, R.,
- 1134 Manjurano, A., Changalucha, J., Elias, J.E., et al. (2017). Seasonal cycling in the gut microbiome of
- 1135 the Hadza hunter-gatherers of Tanzania. Science 357, 802–806.
- 1136 <https://doi.org/10.1126/science.aan4834>.
- 1137 4. Clemente, J.C., Pehrsson, E.C., Blaser, M.J., Sandhu, K., Gao, Z., Wang, B., Magris, M., Hidalgo, G.,
- 1138 Contreras, M., Noya-Alarcón, Ó., et al. (2015). The microbiome of uncontacted Amerindians. Sci.
- 1139 Adv. 1, e1500183–e1500183. <https://doi.org/10.1126/sciadv.1500183>.
- 1140 5. Maghini, D.G., Oduaran, O.H., Olubayo, L.A.I., Cook, J.A., Smyth, N., Mathema, T., Belger, C.W.,
- 1141 Agongo, G., Boua, P.R., Choma, S.S.R., et al. (2025). Expanding the human gut microbiome atlas of
- 1142 Africa. Nature 638, 718–728. <https://doi.org/10.1038/s41586-024-08485-8>.
- 1143 6. Poyet, M., Rühlemann, M., Schaan, A.P., Ma, Y., Moitinho-Silva, L., Wacker, E.M., Jebens, H., Patel,
- 1144 L., Nguyen, L.T.T., Zimmer, A., et al. (2025). Industrialization drives convergent microbial and
- 1145 physiological shifts in the human metaorganism. *Co-submitted along this manuscript*.
- 1146 7. Bobay, L.-M., and Ochman, H. (2018). Factors driving effective population size and pan-genome
- 1147 evolution in bacteria. BMC Evol. Biol., 1–12. <https://doi.org/10.1186/s12862-018-1272-4>.
- 1148 8. Roodgar, M., Good, B.H., Garud, N.R., Martis, S., Avula, M., Zhou, W., Lancaster, S.M., Lee, H.,
- 1149 Babveyh, A., Nesamoney, S., et al. (2021). Longitudinal linked-read sequencing reveals ecological
- 1150 and evolutionary responses of a human gut microbiome during antibiotic treatment. Genome Res 31,
- 1151 1433–1446. <https://doi.org/10.1101/gr.265058.120>.
- 1152 9. Garud, N.R., Good, B.H., Hallatschek, O., and Pollard, K.S. (2019). Evolutionary dynamics of
- 1153 bacteria in the gut microbiome within and across hosts. Plos Biol 17, e3000102.
- 1154 <https://doi.org/10.1371/journal.pbio.3000102>.
- 1155 10. Zhao, S., Lieberman, T.D., Poyet, M., Kauffman, K.M., Gibbons, S.M., Groussin, M., Xavier, R.J., and
- 1156 Alm, E.J. (2019). Adaptive Evolution within Gut Microbiomes of Healthy People. Cell Host Microbe
- 1157 25, 656–667.e8. <https://doi.org/10.1016/j.chom.2019.03.007>.
- 1158 11. Yaffe, E., and Relman, D.A. (2019). Tracking microbial evolution in the human gut using Hi-C reveals
- 1159 extensive horizontal gene transfer, persistence and adaptation. Nat. Microbiol. 16, 472–11.
- 1160 <https://doi.org/10.1038/s41564-019-0625-0>.
- 1161 12. Poyet, M., Groussin, M., Gibbons, S.M., Avila-Pacheco, J., Jiang, X., Kearney, S.M., Perrotta, A.R.,
- 1162 Berdy, B., Zhao, S., Lieberman, T.D., et al. (2019). A library of human gut bacterial isolates paired
- 1163 with longitudinal multiomics data enables mechanistic microbiome research. Nat. Med. 25, 1442–
- 1164 1452. <https://doi.org/10.1038/s41591-019-0559-3>.
- 1165 13. Browne, H.P., Forster, S.C., Anonye, B.O., Kumar, N., Neville, B.A., Stares, M.D., Goulding, D., and
- 1166 Lawley, T.D. (2016). Culturing of “unculturable” human microbiota reveals novel taxa and extensive
- 1167 sporulation. Nature 533, 543–546. <https://doi.org/10.1038/nature17645>.
- 1168 14. Huang, Y., Sheth, R.U., Zhao, S., Cohen, L.A., Dabaghi, K., Moody, T., Sun, Y., Ricaurte, D.,
- 1169 Richardson, M., Velez-Cortes, F., et al. (2023). High-throughput microbial culturomics using
- 1170 automation and machine learning. Nat. Biotechnol., 1–10. [https://doi.org/10.1038/s41587-023-01674-](https://doi.org/10.1038/s41587-023-01674-2)
- 1171 2.
- 1172 15. Forster, S.C., Kumar, N., Anonye, B.O., Almeida, A., Viciani, E., Stares, M.D., Dunn, M., Mkandawire,
- 1173 T.T., Zhu, A., Shao, Y., et al. (2019). A human gut bacterial genome and culture collection for
- 1174 improved metagenomic analyses. Nat. Biotechnol. 37, 186–192. [https://doi.org/10.1038/s41587-018-](https://doi.org/10.1038/s41587-018-0009-7)
- 1175 0009-7.
- 1176 16. Lin, X., Hu, T., Chen, J., Liang, H., Zhou, J., Wu, Z., Ye, C., Jin, X., Xu, X., Zhang, W., et al. (2023).
- 1177 The genomic landscape of reference genomes of cultivated human gut bacteria. Nat. Commun. 14,
- 1178 1663. <https://doi.org/10.1038/s41467-023-37396-x>.

- 1179 17. Goodman, A.L., Kallstrom, G., Faith, J.J., Reyes, A., Moore, A., Dantas, G., and Gordon, J.I. (2011).
1180 Extensive personal human gut microbiota culture collections characterized and manipulated in
1181 gnotobiotic mice. *Proc Natl Acad Sci U S A* 108, 6252–6257. <https://doi.org/10.1073/pnas.1102938108>.
1182 18. Zou, Y., Xue, W., Luo, G., Deng, Z., Qin, P., Guo, R., Sun, H., Xia, Y., Liang, S., Dai, Y., et al. (2019).
1183 1,520 reference genomes from cultivated human gut bacteria enable functional microbiome analyses.
1184 *Nat. Biotechnol.*, 1–15. <https://doi.org/10.1038/s41587-018-0008-8>.
1185 19. Hitch, T.C.A., Masson, J.M., Pauvert, C., Bosch, J., Nüchtern, S., Treichel, N.S., Baloh, M., Razavi,
1186 S., Afrizal, A., Kousetzi, N., et al. (2025). HiBC: a publicly available collection of bacterial strains
1187 isolated from the human gut. *Nat. Commun.* 16, 4203. <https://doi.org/10.1038/s41467-025-59229-9>.
1188 20. Huang, P., Dong, Q., Wang, Y., Tian, Y., Wang, S., Zhang, C., Yu, L., Tian, F., Gao, X., Guo, H., et
1189 al. (2024). Gut microbial genomes with paired isolates from China illustrate probiotic and
1190 cardiometabolic effects. *Cell Genomics*, 100559. <https://doi.org/10.1016/j.xgen.2024.100559>.
1191 21. Groussin, M., Poyet, M., Sistiaga, A., Kearney, S.M., Moniz, K., Noel, M., Hooker, J., Gibbons, S.M.,
1192 Segurel, L., Froment, A., et al. (2021). Elevated rates of horizontal gene transfer in the industrialized
1193 human microbiome. *Cell* 184, 2053-2067.e18. <https://doi.org/10.1016/j.cell.2021.02.052>.
1194 22. Blanco-Míguez, A., Gálvez, E.J.C., Pasolli, E., De Filippis, F., Amend, L., Huang, K.D., Manghi, P.,
1195 Lesker, T.-R., Riedel, T., Cova, L., et al. (2023). Extension of the Segatella copri complex to 13
1196 species with distinct large extrachromosomal elements and associations with host conditions. *Cell*
1197 *Host Microbe* 31, 1804-1819.e9. <https://doi.org/10.1016/j.chom.2023.09.013>.
1198 23. Pasolli, E., Asnicar, F., Manara, S., Zolfo, M., Karcher, N., Armanini, F., Beghini, F., Manghi, P., Tett,
1199 A., Ghensi, P., et al. (2019). Extensive Unexplored Human Microbiome Diversity Revealed by Over
1200 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell*, 1–35.
1201 <https://doi.org/10.1016/j.cell.2019.01.001>.
1202 24. Almeida, A., Nayfach, S., Boland, M., Strozzi, F., Beracochea, M., Shi, Z.J., Pollard, K.S., Sakharova,
1203 E., Parks, D.H., Hugenholtz, P., et al. (2020). A unified catalog of 204,938 reference genomes from
1204 the human gut microbiome. *Nat. Biotechnol.*, 1–25. <https://doi.org/10.1038/s41587-020-0603-3>.
1205 25. Nayfach, S., Shi, Z.J., Seshadri, R., Pollard, K.S., and Kyrpides, N.C. (2019). New insights from
1206 uncultivated genomes of the global human gut microbiome. *Nature* 568, 505–510.
1207 <https://doi.org/10.1038/s41586-019-1058-x>.
1208 26. Karcher, N., Pasolli, E., Asnicar, F., Huang, K.D., Tett, A., Manara, S., Armanini, F., Bain, D.,
1209 Duncan, S.H., Louis, P., et al. (2020). Analysis of 1321 Eubacterium rectale genomes from
1210 metagenomes uncovers complex phylogeographic population structure and subspecies functional
1211 adaptations. *Genome Biol.*, 1–27. <https://doi.org/10.1186/s13059-020-02042-y>.
1212 27. Shaiber, A., and Eren, A.M. (2019). Composite Metagenome-Assembled Genomes Reduce the
1213 Quality of Public Genome Repositories. *mBio* 10, 208–3. <https://doi.org/10.1128/mBio.00725-19>.
1214 28. Meziti, A., Rodriguez-R, L.M., Hatt, J.K., Peña-Gonzalez, A., Levy, K., and Konstantinidis, K.T.
1215 (2021). The Reliability of Metagenome-Assembled Genomes (MAGs) in Representing Natural
1216 Populations: Insights from Comparing MAGs against Isolate Genomes Derived from the Same Fecal
1217 Sample. *Appl. Environ. Microbiol.* 87, e02593-20. <https://doi.org/10.1128/AEM.02593-20>.
1218 29. Mageeney, C.M., Trubl, G., and Williams, K.P. (2022). Improved Mobilome Delineation in
1219 Fragmented Genomes. *Front. Bioinforma.* 2, 866850. <https://doi.org/10.3389/fbinf.2022.866850>.
1220 30. Chen, L.-X., Anantharaman, K., Shaiber, A., Eren, A.M., and Banfield, J.F. (2020). Accurate and
1221 complete genomes from metagenomes. *Genome Res.* 30, 315–333.
1222 <https://doi.org/10.1101/gr.258640.119>.
1223 31. Ramos-Barbero, M.D., Martín-Cuadrado, A.-B., Viver, T., Santos, F., Martínez-García, M., and Antón,
1224 J. (2019). Recovering microbial genomes from metagenomes in hypersaline environments: The
1225 Good, the Bad and the Ugly. *Syst. Appl. Microbiol.* 42, 30–40.
1226 <https://doi.org/10.1016/j.syapm.2018.11.001>.
1227 32. Obregon-Tito, A.J., Tito, R.Y., Metcalf, J., Sankaranarayanan, K., Clemente, J.C., Ursell, L.K., Zech
1228 Xu, Z., Van Treuren, W., Knight, R., Gaffney, P.M., et al. (2015). Subsistence strategies in traditional
1229 societies distinguish gut microbiomes. *Nat. Commun.* 6, 6505. <https://doi.org/10.1038/ncomms7505>.
1230 33. Wibowo, M.C., Yang, Z., Borry, M., Hübner, A., Huang, K.D., Tierney, B.T., Zimmerman, S., Barajas-
1231 Olmos, F., Contreras-Cubas, C., García-Ortiz, H., et al. (2021). Reconstruction of ancient microbial
1232 genomes from the human gut. *Nature* 594, 234–239. <https://doi.org/10.1038/s41586-021-03532-0>.

- 1233 34. Arif, S., Nirmalan, S., Alazizi, A., Mair-Meijers, H., Agyei, A., Afihene, M.Y., Asibey, S.O., Awuku,
1234 Y.A., Duah, A., Plymoth, A., et al. (2025). Host transcriptional responses to gut microbiome variation
1235 arising from urbanism. *Co-submitted along this manuscript*.
- 1236 35. Looft, T., Levine, U.Y., and Stanton, T.B. (2013). Cloacibacillus porcorum sp. nov., a mucin-
1237 degrading bacterium from the swine intestinal tract and emended description of the genus
1238 Cloacibacillus. *Int. J. Syst. Evol. Microbiol.* 63, 1960–1966. <https://doi.org/10.1099/ijs.0.044719-0>.
- 1239 36. Richardson, L., Allen, B., Baldi, G., Beracochea, M., Bileschi, M.L., Burdett, T., Burgin, J., Caballero-
1240 Pérez, J., Cochrane, G., Colwell, L.J., et al. (2023). MGnify: the microbiome sequence data analysis
1241 resource in 2023. *Nucleic Acids Res.* 51, D753–D759. <https://doi.org/10.1093/nar/gkac1080>.
- 1242 37. Chaumeil, P.-A., Mussig, A.J., Hugenholtz, P., and Parks, D.H. (2020). GTDB-Tk: a toolkit to classify
1243 genomes with the Genome Taxonomy Database. *Bioinformatics* 36, 1925–1927.
1244 <https://doi.org/10.1093/bioinformatics/btaz848>.
- 1245 38. Starke, S., Harris, D.M.M., Paulay, A., Aden, K., and Waschina, S. (2025). Comparative analysis of
1246 amino acid auxotrophies and peptidase profiles in non-dysbiotic and dysbiotic small intestinal
1247 microbiomes. *Comput. Struct. Biotechnol. J.* 27, 821–831. <https://doi.org/10.1016/j.csbj.2025.02.004>.
- 1248 39. Mise, K., and Iwasaki, W. (2022). Unexpected absence of ribosomal protein genes from
1249 metagenome-assembled genomes. *ISME Commun.* 2, 118. <https://doi.org/10.1038/s43705-022-00204-6>.
- 1250 40. Rühlemann, M.C., Wacker, E.M., Ellinghaus, D., and Franke, A. (2022). MAGScoT : a fast,
1251 lightweight and accurate bin-refinement tool. *Bioinformatics* 38, 5430–5433.
1252 <https://doi.org/10.1093/bioinformatics/btac694>.
- 1253 41. Smillie, C.S., Smith, M.B., Friedman, J., Cordero, O.X., David, L.A., and Alm, E.J. (2011). Ecology
1254 drives a global network of gene exchange connecting the human microbiome. *Nature* 480, 241–244.
1255 <https://doi.org/10.1038/nature10571>.
- 1256 42. Chen, L., Zhao, N., Cao, J., Liu, X., Xu, J., Ma, Y., Yu, Y., Zhang, X., Zhang, W., Guan, X., et al.
1257 (2022). Short- and long-read metagenomics expand individualized structural variations in gut
1258 microbiomes. *Nat. Commun.* 13, 3175–12. <https://doi.org/10.1038/s41467-022-30857-9>.
- 1259 43. Moss, E.L., Maghini, D.G., and Bhatt, A.S. (2020). Complete, closed bacterial genomes from
1260 microbiomes using nanopore sequencing. *Nat. Biotechnol.*, 1–12. <https://doi.org/10.1038/s41587-020-0422-6>.
- 1261 44. Morel, B., Williams, T.A., Stamatakis, A., and Szöllösi, G.J. (2024). AleRax: a tool for gene and
1262 species tree co-estimation and reconciliation under a probabilistic model of gene duplication, transfer,
1263 and loss. *Bioinformatics* 40, btae162. <https://doi.org/10.1093/bioinformatics/btae162>.
- 1264 45. Mainprize, I.L., Bean, J.D., Bouwman, C., Kimber, M.S., and Whitfield, C. (2013). The UDP-glucose
1265 Dehydrogenase of *Escherichia coli* K-12 Displays Substrate Inhibition by NAD That Is Relieved by
1266 Nucleotide Triphosphates. *J. Biol. Chem.* 288, 23064–23074.
1267 <https://doi.org/10.1074/jbc.M113.486613>.
- 1268 46. Lehrer, J., Vigeant, K.A., Tatar, L.D., and Valvano, M.A. (2007). Functional Characterization and
1269 Membrane Topology of *Escherichia coli* WecA, a Sugar-Phosphate Transferase Initiating the
1270 Biosynthesis of Enterobacterial Common Antigen and O-Antigen Lipopolysaccharide. *J. Bacteriol.*
1271 189, 2618–2628. <https://doi.org/10.1128/JB.01905-06>.
- 1272 47. Virolle, C., Goldlust, K., Djermoun, S., Bigot, S., and Lesterlin, C. (2020). Plasmid Transfer by
1273 Conjugation in Gram-Negative Bacteria: From the Cellular to the Community Level. *Genes* 11, 1239.
1274 <https://doi.org/10.3390/genes11111239>.
- 1275 48. Wolff, R., and Garud, N.R. Pervasive selective sweeps across human gut microbiomes.
1276
- 1277 49. Groussin, M., Poyet, M., Sistiaga, A., Kearney, S.M., Moniz, K., Noel, M., Hooker, J., Gibbons, S.M.,
1278 Ségurel, L., Froment, A., et al. (2021). Elevated rates of horizontal gene transfer in the industrialized
1279 human microbiome. *Cell* 184, 2053-2067.e18. <https://doi.org/10.1016/j.cell.2021.02.052>.
- 1280 50. Haines-Menges, B., Whitaker, W.B., and Boyd, E.F. (2014). Alternative Sigma Factor RpoE Is
1281 Important for *Vibrio parahaemolyticus* Cell Envelope Stress Response and Intestinal Colonization.
1282 *Infect. Immun.* 82, 3667–3677. <https://doi.org/10.1128/IAI.01854-14>.
- 1283 51. Gutierrez, J., Smith, R., and Pogliano, K. (2010). SpoIID-Mediated Peptidoglycan Degradation Is
1284 Required throughout Engulfment during *Bacillus subtilis* Sporulation. *J. Bacteriol.* 192, 3174–3186.
1285 <https://doi.org/10.1128/JB.00127-10>.
- 1286 52. Heimesaat, M.M., Schmidt, A.-M., Mousavi, S., Escher, U., Tegtmeyer, N., Wessler, S., Gadermaier,
1287 G., Briza, P., Hofreuter, D., Bereswill, S., et al. (2020). Peptidase PepP is a novel virulence factor of
1288

- 1289 *Campylobacter jejuni* contributing to murine campylobacteriosis. *Gut Microbes* 12, 1770017.
1290 <https://doi.org/10.1080/19490976.2020.1770017>.
- 1291 53. Paterson, G.K., Cone, D.B., Northen, H., Peters, S.E., and Maskell, D.J. (2009). Deletion of the gene
1292 encoding the glycolytic enzyme triosephosphate isomerase (*tpi*) alters morphology of *Salmonella*
1293 *enterica* serovar Typhimurium and decreases fitness in mice. *FEMS Microbiol. Lett.* 294, 45–51.
1294 <https://doi.org/10.1111/j.1574-6968.2009.01553.x>.
- 1295 54. Xia, Y., Wang, D., Pan, X., Xia, B., Weng, Y., Long, Y., Ren, H., Zhou, J., Jin, Y., Bai, F., et al.
1296 (2020). TpiA is a Key Metabolic Enzyme That Affects Virulence and Resistance to Aminoglycoside
1297 Antibiotics through CrcZ in *Pseudomonas aeruginosa*. *mBio* 11, e02079-19.
1298 <https://doi.org/10.1128/mBio.02079-19>.
- 1299 55. Urbonavičius, J., Durand, J.M.B., and Björk, G.R. (2002). Three Modifications in the D and T Arms of
1300 tRNA Influence Translation in *Escherichia coli* and Expression of Virulence Genes in *Shigella flexneri*.
1301 *J. Bacteriol.* 184, 5348–5357. <https://doi.org/10.1128/JB.184.19.5348-5357.2002>.
- 1302 56. Noel, H.R., Keerthi, S., Ren, X., Winkelman, J.D., Troutman, J.M., and Palmer, L.D. (2024). Genetic
1303 synergy between *Acinetobacter baumannii* undecaprenyl phosphate biosynthesis and the Mla system
1304 impacts cell envelope and antimicrobial resistance. *mBio* 15, e02804-23.
1305 <https://doi.org/10.1128/mbio.02804-23>.
- 1306 57. Verstraete, M.M., Perez-Borraero, C., Brown, K.L., Heinrichs, D.E., and Murphy, M.E.P. (2018). SbnI
1307 is a free serine kinase that generates -phospho-l-serine for staphyloferrin B biosynthesis in. *J. Biol.*
1308 *Chem.* 293, 6147–6160. <https://doi.org/10.1074/jbc.RA118.001875>.
- 1309 58. Révora, V., Marchesini, M.I., and Comerci, D.J. (2020). *Brucella abortus* Depends on L -Serine
1310 Biosynthesis for Intracellular Proliferation. *Infect. Immun.* 88, e00840-19.
1311 <https://doi.org/10.1128/IAI.00840-19>.
- 1312 59. Croucher, N.J., Page, A.J., Connor, T.R., Delaney, A.J., Keane, J.A., Bentley, S.D., Parkhill, J., and
1313 Harris, S.R. (2015). Rapid phylogenetic analysis of large samples of recombinant bacterial whole
1314 genome sequences using Gubbins. *Nucleic Acids Res* 43, e15–e15.
1315 <https://doi.org/10.1093/nar/gku1196>.
- 1316 60. Sonnenburg, J.L., and Sonnenburg, E.D. (2019). Vulnerability of the industrialized microbiota.
1317 *Science* 366, eaaw9255. <https://doi.org/10.1126/science.aaw9255>.
- 1318 61. Vatanen, T., Kostic, A.D., d’Hennezel, E., Siljander, H., Franzosa, E.A., Yassour, M., Kolde, R.,
1319 Vlamakis, H., Arthur, T.D., Hämäläinen, A.-M., et al. (2016). Variation in Microbiome LPS
1320 Immunogenicity Contributes to Autoimmunity in Humans. *Cell* 165, 842–853.
1321 <https://doi.org/10.1016/j.cell.2016.04.007>.
- 1322 62. Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads.
1323 *EMBnet J*, 10–12.
- 1324 63. Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina
1325 sequence data. *Bioinformatics* 30, 2114–2120. <https://doi.org/10.1093/bioinformatics/btu170>.
- 1326 64. Bankevich, A., Nurk, S., Antipov, D., Gurevich, A.A., Dvorkin, M., Kulikov, A.S., Lesin, V.M.,
1327 Nikolenko, S.I., Pham, S., Prjibelski, A.D., et al. (2012). SPAdes: a new genome assembly algorithm
1328 and its applications to single-cell sequencing. *J Comput Biol* 19, 455–477.
1329 <https://doi.org/10.1089/cmb.2012.0021>.
- 1330 65. Boetzer, M., Henkel, C.V., Jansen, H.J., Butler, D., and Pirovano, W. (2011). Scaffolding pre-
1331 assembled contigs using SSPACE. *Bioinformatics* 27, 578–579.
1332 <https://doi.org/10.1093/bioinformatics/btq683>.
- 1333 66. Nadalin, F., Vezzi, F., and Policriti, A. (2012). GapFiller: a de novo assembly approach to fill the gap
1334 within paired reads. *BMC Bioinformatics* 13, S8. <https://doi.org/10.1186/1471-2105-13-S14-S8>.
- 1335 67. Li, D., Liu, C.-M., Luo, R., Sadakane, K., and Lam, T.-W. (2015). MEGAHIT: an ultra-fast single-node
1336 solution for large and complex metagenomics assembly via succinct *de Bruijn* graph. *Bioinformatics*
1337 31, 1674–1676. <https://doi.org/10.1093/bioinformatics/btv033>.
- 1338 68. Wu, Y.-W., Simmons, B.A., and Singer, S.W. (2016). MaxBin 2.0: an automated binning algorithm to
1339 recover genomes from multiple metagenomic datasets. *Bioinformatics* 32, 605–607.
1340 <https://doi.org/10.1093/bioinformatics/btv638>.
- 1341 69. Kang, D.D., Li, F., Kirton, E., Thomas, A., Egan, R., An, H., and Wang, Z. (2019). MetaBAT 2: an
1342 adaptive binning algorithm for robust and efficient genome reconstruction from metagenome
1343 assemblies. *PeerJ* 7, e7359. <https://doi.org/10.7717/peerj.7359>.

- 1344 70. Alneberg, J., Bjarnason, B.S., De Bruijn, I., Schirmer, M., Quick, J., Ijaz, U.Z., Lahti, L., Loman, N.J.,
1345 Andersson, A.F., and Quince, C. (2014). Binning metagenomic contigs by coverage and composition.
1346 *Nat. Methods* 11, 1144–1146. <https://doi.org/10.1038/nmeth.3103>.
- 1347 71. Nissen, J.N., Johansen, J., Allesøe, R.L., Sønderby, C.K., Armenteros, J.J.A., Grønbech, C.H.,
1348 Jensen, L.J., Nielsen, H.B., Petersen, T.N., Winther, O., et al. (2021). Improved metagenome binning
1349 and assembly using deep variational autoencoders. *Nat. Biotechnol.* 39, 555–560.
1350 <https://doi.org/10.1038/s41587-020-00777-4>.
- 1351 72. Parks, D.H., Imelfort, M., Skennerton, C.T., Hugenholtz, P., and Tyson, G.W. (2015). CheckM:
1352 assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes.
1353 *Genome Res* 25, 1043–1055. <https://doi.org/10.1101/gr.186072.114>.
- 1354 73. Srivastava, A., Malik, L., Smith, T., Sudbery, I., and Patro, R. (2019). Alevin efficiently estimates
1355 accurate gene abundances from dscRNA-seq data. *Genome Biol.* 20, 65.
1356 <https://doi.org/10.1186/s13059-019-1670-y>.
- 1357 74. Zimmermann, J., Kaleta, C., and Waschina, S. (2021). gapseq: informed prediction of bacterial
1358 metabolic pathways and reconstruction of accurate metabolic models. *Genome Biol* 22, 81–35.
1359 <https://doi.org/10.1186/s13059-021-02295-1>.
- 1360 75. Starke, S., Harris, D.M.M., Zimmermann, J., Schuchardt, S., Oumari, M., Frank, D., Bang, C.,
1361 Rosenstiel, P., Schreiber, S., Frey, N., et al. (2023). Amino acid auxotrophies in human gut bacteria
1362 are linked to higher microbiome diversity and long-term stability. *ISME J.* 17, 2370–2380.
1363 <https://doi.org/10.1038/s41396-023-01537-3>.
- 1364 76. Feldgarden, M., Brover, V., Gonzalez-Escalona, N., Frye, J.G., Haendiges, J., Haft, D.H., Hoffmann,
1365 M., Pettengill, J.B., Prasad, A.B., Tillman, G.E., et al. (2021). AMRFinderPlus and the Reference
1366 Gene Catalog facilitate examination of the genomic links among antimicrobial resistance, stress
1367 response, and virulence. *Sci. Rep.* 11, 12728. <https://doi.org/10.1038/s41598-021-91456-0>.
- 1368 77. Liu, B., Zheng, D., Jin, Q., Chen, L., and Yang, J. (2019). VFDB 2019: a comparative pathogenomic
1369 platform with an interactive web interface. *Nucleic Acids Res.* 47, D687–D692.
1370 <https://doi.org/10.1093/nar/gky1080>.
- 1371 78. Zheng, J., Ge, Q., Yan, Y., Zhang, X., Huang, L., and Yin, Y. (2023). dbCAN3: automated
1372 carbohydrate-active enzyme and substrate annotation. *Nucleic Acids Res.* 51, W115–W121.
1373 <https://doi.org/10.1093/nar/gkad328>.
- 1374 79. Abby, S.S., Denise, R., and Rocha, E.P.C. (2024). Identification of Protein Secretion Systems in
1375 Bacterial Genomes Using MacSyFinder Version 2. In *Bacterial Secretion Systems Methods in*
1376 *Molecular Biology.*, L. Journet and E. Cascales, eds. (Springer US), pp. 1–25.
1377 https://doi.org/10.1007/978-1-0716-3445-5_1.
- 1378 80. Xie, Z., and Tang, H. (2017). ISEScan: automated identification of insertion sequence elements in
1379 prokaryotic genomes. *Bioinformatics* 33, 3340–3347. <https://doi.org/10.1093/bioinformatics/btx433>.
- 1380 81. Camargo, A.P., Roux, S., Schulz, F., Babinski, M., Xu, Y., Hu, B., Chain, P.S.G., Nayfach, S., and
1381 Kyrpides, N.C. (2023). Identification of mobile genetic elements with geNomad. *Nat. Biotechnol.*
1382 <https://doi.org/10.1038/s41587-023-01953-y>.
- 1383 82. Camarillo-Guerrero, L.F., Almeida, A., Rangel-Pineros, G., Finn, R.D., and Lawley, T.D. (2021).
1384 Massive expansion of human gut bacteriophage diversity. *Cell* 184, 1098–1109.e9.
1385 <https://doi.org/10.1016/j.cell.2021.01.029>.
- 1386 83. Nayfach, S., Camargo, A.P., Schulz, F., Eloë-Fadrosch, E., Roux, S., and Kyrpides, N.C. (2021).
1387 CheckV assesses the quality and completeness of metagenome-assembled viral genomes. *Nat.*
1388 *Biotechnol.* 39, 578–585. <https://doi.org/10.1038/s41587-020-00774-7>.
- 1389 84. Preska Steinberg, A., and Kussell, E. (2025). How recombination and clonal evolution shape bacterial
1390 lineages and genomes. *GENETICS*, iyaf115. <https://doi.org/10.1093/genetics/iyaf115>.
- 1391 85. Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X., and Ruden,
1392 D.M. (2012). A program for annotating and predicting the effects of single nucleotide polymorphisms,
1393 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* 6,
1394 80–92. <https://doi.org/10.4161/fly.19695>.
- 1395 86. Minh, B.Q., Schmidt, H.A., Chernomor, O., Schrempf, D., Woodhams, M.D., Von Haeseler, A., and
1396 Lanfear, R. (2020). IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the
1397 Genomic Era. *Mol. Biol. Evol.* 37, 1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
- 1398 87. Seemann, T. (2014). Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30, 2068–2069.
1399 <https://doi.org/10.1093/bioinformatics/btu153>.

- 1400 88. Steinegger, M., and Söding, J. (2017). MMseqs2 enables sensitive protein sequence searching for
1401 the analysis of massive data sets. *Nat. Biotechnol.* 35, 1026–1028. <https://doi.org/10.1038/nbt.3988>.
- 1402 89. Cantalapiedra, C.P., Hernández-Plaza, A., Letunic, I., Bork, P., and Huerta-Cepas, J. (2021).
1403 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain Prediction at the
1404 Metagenomic Scale. *Mol. Biol. Evol.* 38, 5825–5829. <https://doi.org/10.1093/molbev/msab293>.
- 1405 90. Katoh, K., and Standley, D.M. (2013). MAFFT Multiple Sequence Alignment Software Version 7:
1406 Improvements in Performance and Usability. *Mol. Biol. Evol.* 30, 772–780.
1407 <https://doi.org/10.1093/molbev/mst010>.
- 1408 91. Kislyuk, A.O., Haegeman, B., Bergman, N.H., and Weitz, J.S. (2011). Genomic fluidity: an integrative
1409 view of gene diversity within microbial populations. *BMC Genomics* 12, 32.
1410 <https://doi.org/10.1186/1471-2164-12-32>.
- 1411 92. Dewar, A.E., Hao, C., Belcher, L.J., Ghouli, M., and West, S.A. (2024). Bacterial lifestyle shapes
1412 pangenomes. *Proc. Natl. Acad. Sci.* 121, e2320170121. <https://doi.org/10.1073/pnas.2320170121>.
- 1413 93. Szöllösi, G.J., Rosikiewicz, W., Boussau, B., Tannier, E., and Daubin, V. (2013). Efficient exploration
1414 of the space of reconciled gene trees. *Syst Biol* 62, 901–912. <https://doi.org/10.1093/sysbio/syt054>.
- 1415 94. Huerta-Cepas, J., Serra, F., and Bork, P. (2016). ETE 3: Reconstruction, Analysis, and Visualization
1416 of Phylogenomic Data. *Mol. Biol. Evol.* 33, 1635–1638. <https://doi.org/10.1093/molbev/msw046>.
- 1417 95. Hoang, D.T., Chernomor, O., Von Haeseler, A., Minh, B.Q., and Vinh, L.S. (2018). UFBoot2:
1418 Improving the Ultrafast Bootstrap Approximation. *Mol. Biol. Evol.* 35, 518–522.
1419 <https://doi.org/10.1093/molbev/msx281>.
- 1420 96. Kluge, A.G., and Farris, J.S. (1969). Quantitative Phyletics and the Evolution of Anurans. *Syst. Biol.*
1421 18, 1–32. <https://doi.org/10.1093/sysbio/18.1.1>.
- 1422 97. Tung Ho, L.S., and Ané, C. (2014). A Linear-Time Algorithm for Gaussian and Non-Gaussian Trait
1423 Evolution Models. *Syst. Biol.* 63, 397–408. <https://doi.org/10.1093/sysbio/syu005>.

1424

1425 **List of Supplementary materials**

1426

1427 **Supplementary Figures**

- 1428 ● Supp. Figure 1 – Distribution and functional predictions of GMbC genomes reveal host
1429 and strain-level diversity.
- 1430 ● Supp. Figure 2 – Pairwise comparisons of genomic features between MAG and isolate
1431 genomes sampled from the same donor highlight the recovery limitations of MAGs.
- 1432 ● Supp. Figure 3 – Pangenome characteristics of the ten species studied for convergent
1433 adaptation to industrialization.
- 1434 ● Supp. Figure 4 – Gene tree-Species tree reconciliations reveal increasing gene content
1435 along lineages of industrialized hosts.
- 1436 ● Supp. Figure 5 – Gene enrichment analysis based on industrialization level for BTHE,
1437 BXYL, PDOR, PMER and BWEX.
- 1438 ● Supp. Figure 6 – Convergent elevation of gene-level Ka/Ks based on industrialization
1439 status.
- 1440 ● Supp. Figure 7 – Convergent signals of SNV-host lifestyle associations across bacterial
1441 species at the level of KEGG KO categories.

1442

1443 **Supplementary Tables**

- 1444 ● Supp. Table 1 – Metadata of isolate genomes and donor participants.
- 1445 ● Supp. Table 2 – Profiles and comparisons of genomic features between paired MAGs and
1446 isolate genomes.
- 1447 ● Supp. Table 3 – Pangenome fluidity estimates across ten species.
- 1448 ● Supp. Table 4 – Counts of Gene tree - Species tree reconciliation events (speciations,
1449 duplications, losses, transfers, presence, originations, copies and singletons) for the ten
1450 species.
- 1451 ● Supp. Table 5 – Statistics and results of gene enrichment analyses performed on 90%
1452 similarity gene families in the ten species.
- 1453 ● Supp. Table 6 – Gene-level Ka/Ks estimates of 90% similarity gene families for the ten
1454 species.
- 1455 ● Supp. Table 7 – Statistics and results of variant-level associations with host
1456 industrialization for BFRA, BOVA, BUNI, PDIS and PVUL.

1457

1458

1459

1460 **Acknowledgments**

1461 This work was supported by grants from the Center for Microbiome Informatics and Therapeutics
1462 at MIT and the Rasmussen Family Foundation.

1463 M.P, M.R and M.G. received support from the Deutsche Forschungsgemeinschaft (DFG - German
1464 Science Foundation) within the Collaborative Research Center (CRC) 1182 on “The Origin and
1465 Function of Metaorganisms” (Project-ID 261376515 – SFB 1182, project C5.1 to M.G., project
1466 C5.2 to M.P., project C5.3 to M.R.).

1467 The study received infrastructure support from the DFG (German Science Foundation) within the
1468 Cluster of Excellence 2167 “PrecisionMedicine in Chronic Inflammation (PMI)” (EXC 2167-
1469 390884018).

1470 M.G. received funding from the European Research Council (ERC) under the European Union’s
1471 Horizon 2020 research and innovation programme (CoG VESICULOME, Grant agreement No.
1472 101126254).

1473 M.G., M.P. and A.F. received funding from the DFG – Project number 426660215 with the RU
1474 5042 miTarget on “The microbiome as a therapeutic target in inflammatory bowel disease”,
1475 subprojects TP01 “Targeting intestinal yeasts and pathogenic yeast-responsive CD4+ T cells in
1476 Crohn’s disease” and TP02 “Ecology and function of synthetic bacterial communities for the
1477 understanding and modulation of IBD-associated microbiomes”.

1478 The project received funding from the European Union under the Horizon Europe grant agreement
1479 No. 101095470 (project miGut-Health). Views and opinions expressed are however those of the
1480 author(s) only and do not necessarily reflect those of the European Union nor European Health
1481 and Digital Executive Agency (HaDEA). Neither the European Union nor HaDEA can be held
1482 responsible for them.

1483 G.J.Sz and L.L.Sz. are grateful for the help and support provided by the Scientific Computing and
1484 Data Analysis section of Core Facilities at OIST.

1485 R. J. X. acknowledges funding from the National Institutes of Health (NIH) (Project Center for the
1486 Study of Inflammatory Bowel Disease at Massachusetts General Hospital - DK043351)

1487 We thank Tamara Mason and the team at the Walkup Sequencing platform at the Broad Institute
1488 for support on sequencing efforts.

1489

1490 **Declaration of interests**

1491 R.J.X. is a co-founder of Convergence Bio, board director at MoonLake Immunotherapeutics, a
1492 consultant to Nestlé, and a member of the advisory boards at MagnetBiomedicine and Arena
1493 Bioworks. No organizations listed above provided funding for this study.

1494

1495 **Author contributions**

1496 Designed this study: M.P, M.G, M.R, E.J.A

1497 Field administrative work & collection of data and samples: M.P, M.G, V.J, A.F, A.A, M.Y.A, S.O.A,
1498 Y.A.A, A.D, Y.A.N, F.I, Y.L.A.L, T.M.P, C.O, J.R, I.E.M

1499 Performed processing of biospecimens and molecular work: M.P

1500 Biosample and Data curation: M.P, M.G, M.R, L.M, H.J, J.C

1501 Data analysis: M.R, M.G, M.P, E.J.A, L.L.S, S.W, L.M-S, L.K.M, J.F.C, J.B, A.F, G.J.S

1502 Supervision: M.G, M.P, E.J.A

1503 Funding acquisition: M.P, M.G, E.J.A, R.J.X, A.F, J.B, G.J.S
1504 Writing, original draft: M.G, M.P, M.R
1505 Writing, review & editing: all authors
1506
1507
1508
1509
1510

1511 **Figure 1 – Geographic, species, strain, and host lifestyle diversity in the GMbC collection**
1512 **of human gut bacterial isolate genomes**

1513

1514 H. Overview of sampling, preservation, culturing, isolation, and sequencing procedures for
1515 gut bacterial genomes (see Methods).

1516 I. Lifestyle and microbiome diversity of donors used for culturing and isolating gut bacteria
1517 in the context of the broader GMbC + BIO-ML cohort. Top panel: dimensional reduction
1518 analysis of various lifestyle factors (see Methods). Donors used for culturing are shown in
1519 larger symbols with dark border. GMbC donors are shown in circles, BIO-ML donors are
1520 shown in triangles. Spearman correlations between the first two PCs and individual
1521 lifestyle factors are shown on the right. Alpha diversity (measured with Faith PD index)
1522 and beta diversity (unweighted UniFrac) of GMbC and BIO-ML isolate donors and
1523 participants are shown in bottom panels.

1524 J. Phylogenomic tree of representative genomes from 434 species-level genome bins
1525 (SGBs). Inner ring shows overlap with external genome collections (UHGG v2, GMbC
1526 MAGs, BIO-ML). Middle ring indicates host lifestyle origin (industrialized or non-
1527 industrialized). Outer ring shows country distribution and isolate genome counts per SGB.
1528 Clade colors represent phyla.

1529 K. Isolate genome, strain bin, and SGB counts by country and host lifestyle. Strain bins group
1530 genomes from the same donor with >99% similarity (see Methods). Counts that include
1531 isolate genomes of the BIO-ML collection per host lifestyle are also shown. BIO-ML isolate
1532 genomes were generated following the same pipeline as described in A (see Methods).

1533 L. Distribution of strain bin counts across SGBs, localities and individual hosts. Colors denote
1534 country.

1535 M. Ten bacterial species with ≥ 8 strain bins sampled from industrialized or non-industrialized
1536 hosts. Barplots show isolate and strain bin counts per lifestyle.

1537 N. Phylogenetic trees of representative strain bin genomes for the 10 species in panel E. Tip
1538 points indicate host lifestyle; labels show country/locality and are color-coded by country.
1539 Trees are midpoint-rooted. Branch length scales are in expected number of substitutions
1540 per site.

1541

1542 **Figure 2 – Isolate genomes recover more genomic features and HGT events than MAGs.**

1543

1544 E. Number and quality scores of MAG–isolate genome pairs. Pairs originate from the same
1545 donor sample and species.

1546 F. Heatmap comparing genomic feature counts across all genome pairs. Genera and species
1547 of genome pairs are shown on the left and right side of the heatmap, respectively. Genomic
1548 features are shown in columns, and are grouped in four categories: genomic size,
1549 metabolism, key functions and mobile genetic elements (MGEs) & machineries. For each
1550 pair, the difference in counts between the isolate genome and the MAG was calculated
1551 and normalized to the count in the isolate genome. Features with higher counts in the
1552 isolate genome or in the MAG are shown along a gradient of red to blue, respectively.

1553 G. Summary statistics of feature differences across all pairs.

1554 H. Comparison of HGT events. Between-species horizontal gene transfers (HGTs) were
1555 detected across isolate genomes, and across MAGs separately (see Methods). Genomes
1556 of MAG-isolate genome pairs cluster in 87 SGBs. Ratio of species pairs with detected
1557 HGTs ($n \geq 1$ HGT) were compared with a proportion test (Two proportion Z-test, ***: $p =$
1558 $7.96e-33$). Edges in the network indicate that at least 1 HGT was detected between
1559 species (nodes).

1560

1561

1562 **Figure 3 – Industrialized host strains exhibit larger proteomes and signatures of relaxed**
1563 **selection**

1564

1565 C. Comparison of proteome size (coding gene counts) between strains of host with
1566 industrialized vs. non-industrialized lifestyles (in purple and green, respectively) across
1567 the 10 species presented in Figure 1. Counts were statistically compared while accounting
1568 for phylogeny (phyloglm function, see Methods) (***: p-value < 0.001; **: p-value < 0.01;
1569 *: p-value < 0.05; NS: non-significant – this legend applies to all other panels). P-values
1570 were combined with the Fisher’s method to test for cross-species evidence of differences
1571 in proteome size against the null hypothesis. This p-value is shown on the right of the
1572 panel.

1573 D. Comparison of pangenome fluidity among industrialization- and non-industrialization-
1574 associated strains. The ratio of shared genes was calculated for strain bin pairs, using
1575 representative genomes. P-values were combined with the Fisher’s method (p-value
1576 shown on the right of the panel)

1577

1578

1579 **Figure 4 – Recent and HGT-driven gene gains promote proteome expansion in**
1580 **industrialized strains**

1581

1582 D. Species tree - gene tree reconciliations were sampled to detect and count per-branch
1583 events of gene transfer, loss, origination and speciation (see Methods). Counts of per-
1584 branch gene loss (left column) and gain (middle), and differences between gain and loss
1585 counts (right) were compared between host lifestyle categories (industrialized: purple
1586 area; non-industrialized: green area). Gene gains were defined as the sum of gene
1587 transfer and origination events. Top row: counts aggregated across all branches. Middle
1588 row: Counts of internal branches. Bottom row: counts of terminal (tip) branches. For each
1589 species, median counts are shown, with intervals ranging from the 25th to the 75th
1590 quantiles. Plain points indicate species for which the difference in counts between host
1591 lifestyle categories is significantly different (Wilcoxon tests). Species are colored-coded.

1592 E. Correlation between per-branch gain–loss differences and HGT counts, broken down by
1593 internal (green) and terminal (orange) branches. All correlations are statistically significant
1594 ($p\text{-val} < 0.001$; Spearman correlation tests).

1595 F. Increasing gene content along lineages of industrialized hosts. The panel depicts the
1596 evolution of the number of genes along the phylogeny of BOVA, BTHE, BXYL, PDOR and
1597 PVUL, based on the reconciliation-aware reconstruction of ancestral gene contents.
1598 Reconciliations were calculated from the set of gene families present in at least four StGBs
1599 (tips of the tree). These 5 species have significant differences in proteome size (Figure 3)
1600 and in gene gains along terminal branches between host lifestyles. Data for the other five
1601 species, which show similar trends, is presented in Supp. Fig. 4.

1602

1603

1604 **Figure 5 – Genes differentially enriched between host industrialized and non-industrialized**
1605 **lifestyles**

1606

1607 C. Gene enrichment analysis based on categorical and continuous levels of industrialization.

1608 Gene profiles were coded as presence/absence data and were correlated to host
1609 industrialization status encoded as a binary variable. Differential enrichment was tested
1610 while controlling for phylogeny. Significant hits (q-value < 0.05) are colored based
1611 on host lifestyle (purple: industrialized; green: non-industrialized). Non-significant genes
1612 are colored in grey. Top row: volcano plots showing all genes. The number of statistically
1613 significant genes is shown next to each species acronym. Significantly differentially
1614 enriched genes validated by measuring correlations with PC1 Lifestyle rather than
1615 industrialization status as a binary variable are shown in plain circles. Significant genes
1616 not validated with PC1 Lifestyle are shown as empty circles. Bottom row: top 10 most
1617 differentially enriched genes, for each lifestyle category. Genes with similar
1618 presence/absence profiles are collapsed into a single gene cluster. Gene labels indicate
1619 the number (n) of 90% gene families collapsed together.

1620 Data for 5 species are shown. These species harbor most of the significant hits. Data for
1621 the other 5 species is shown in Supp. Fig. 5.

1622 D. Gene families (90% and 50% similarity gene clusters on top and bottom panels,
1623 respectively) with signals of differential enrichment across multiple species. Most gene
1624 families show convergent signals of differential enrichment based on host lifestyle
1625 (enrichment in one of two lifestyle categories across multiple species).

1626

1627

1628 **Figure 6 – Lifestyle-specific signals of positive selection at the gene level**

1629

1630 G. Gene-level Ka/Ks values across species and host lifestyles (90% similarity gene families).

1631 For each lifestyle category and each gene, Ka/Ks values were computed for all pairs of
1632 codon-aligned gene sequences. Median Ka/Ks values are reported.

1633 H. Distribution of species-level median Ka/Ks values, aggregated across all genes.

1634 I. Counts of genes with median Ka/Ks values ≥ 1 (positive selection) across host lifestyle
1635 categories.

1636 J. Percentage of genes with median Ka/Ks values ≥ 1 across host lifestyle categories.

1637 K. Operons, 50% similarity gene clusters and KEGG KOs with convergent signals of positive
1638 selection across 2+ species.

1639 L. Top genes with highest absolute differences in median Ka/Ks values between
1640 industrialized and non-industrialized lifestyle categories.

1641

1642

1643 **Figure 7 – Patterns of single nucleotide polymorphisms reveal genes and operons that**
1644 **undergo cross-species parallel evolution associated with host lifestyle**

1645

1646 C. Associations between single nucleotide variants (SNVs) and host lifestyle categories were
1647 calculated while accounting for phylogeny (see Methods). Associations were calculated
1648 for the 5 species with enough genome sample size to yield sufficient statistical power. Hits
1649 (q-value < 0.05) were cross-validated using GMbC shotgun metagenomes (n = 1,015, see
1650 Methods) that were sampled from diverse geographies worldwide, including those from
1651 which isolate genomes originate (reference). Hits validated by metagenomes are shown
1652 as diamonds and annotated.

1653 D. Convergent signals of SNV-host lifestyle associations across bacterial species. Each tile
1654 represents an operon–species association and contains the names of genes within that
1655 operon that contain host-lifestyle associated SNVs. Operons are shown along the x-axis,
1656 and species along the y-axis, cells are color-coded by species identity. Black points show
1657 operons in which similar genes contain host-lifestyle associated SNVs across species. All
1658 SNV data can be found in Supp. Table 7.

1659

1660

1661

1662 **Supplementary Figure 1 – Distribution and functional predictions of GMbC genomes reveal**
1663 **host and strain-level diversity.**

- 1664 A. Distribution of SGB and genome counts across localities and individual hosts. Colors
1665 denote country.
- 1666 B. Phylogenomic tree of representative genomes from 434 species-level genome bins
1667 (SGBs). Outer rings show annotations of functional and phenotypic categories predicted
1668 by Traitar. For each trait and SGB, trait prevalence was first calculated across all genomes
1669 of the corresponding strain-level genome bins (StGBs), and SGB-level conservation
1670 scores were then obtained by averaging StGB prevalence values. Trait conservation
1671 scores are visualized using an opacity gradient ('score' variable).
- 1672 C. Heatmap of strain-level amino acid auxotrophy variability across species (y axis). Values
1673 show the variance in auxotrophy predictions across genomes of a given species. Only
1674 genomes with quality score higher than 95% were included in the variance calculations.
1675 Auxotrophy predictions were performed from gapseq-derived genome-scale metabolic
1676 models (see Methods).

1678 **Supplementary Figure 2 – Pairwise comparisons of genomic features between MAG and**
1679 **isolate genomes sampled from the same donor highlight the recovery limitations of MAGs.**

- 1680 A. Metagenomic abundance of species included in the MAG vs. isolate genome comparison.
1681 Each data point represents a pair of MAG and isolate genome from the same species,
1682 sampled within the same host. Species are colored by phylum.
- 1683 B. Counts of various genomic features in the paired MAG and isolate genomes (n = 147).
1684 Genome pairs (species shown along the x axis) are grouped by genus. Empty green
1685 circles depict MAG estimates, black filled circles depict isolate genome estimates. See
1686 Methods for a full description of the functional categories and their genomic profiling.

1688 **Supplementary Figure 3 – Pangenome characteristics of the ten species studied for**
1689 **convergent adaptation to industrialization.**

- 1690 A. Count of the number of genes in the core-, accessory- and cloud-genome of the ten
1691 species of interest. See Methods for the definition and reconstruction of each pangenome
1692 category.
- 1693 B. Count of the number of genes in pangenome categories as a function of the number of
1694 strain-level genome bins (StGBs) per species. Linear regressions per pangenome
1695 categories were calculated – core-genome: beta = -2.8, p-val = 0.8; accessory-genome:
1696 beta = 4.8, p-val = 0.9; cloud-genome: beta = 282, p-val = 5.59e-05. Size of core and
1697 accessory genomes are stable across species irrespective of the number of StGBs. The
1698 size of the cloud genome is positively correlated to the number of sampled StGBs.

1700 **Supplementary Figure 4 – Gene tree-Species tree reconciliations reveal increasing gene**
1701 **content along lineages of industrialized hosts.**

1702 The panel depicts the evolution of the number of genes along the phylogeny of BFRA, BUNI,
1703 PDIS, PMER and BWEX, based on the reconciliation-aware reconstruction of ancestral gene
1704 contents. Reconciliations were calculated from the set of gene families present in at least four
1705 StGBs (tips of the tree). Overall trends for an increase in gene content along lineages occurring

1706 in industrialized hosts is observed (Fig. 4), and is statistically significant for BFRA and PMER (Fig.
1707 4).

1708

1709

1710 **Supplementary Figure 5 – Gene enrichment analysis based on industrialization level for**
1711 **BTHE, BXYL, PDOR, PMER and BWEX.**

1712 Gene enrichment analysis performed for BTHE, BXYL, PDOR, PMER and BWEX, based on
1713 categorical and continuous levels of industrialization. Data is presented as in Fig. 5. In brief, gene
1714 profiles were coded as presence/absence data and were correlated to host industrialization
1715 status, controlling for phylogeny. Significant hits (q-value < 0.05) are colored coded based on host
1716 lifestyle. Volcano plots showing data for all genes (top) and effect size plots of top hits (bottom)
1717 are shown. The number of statistically significant genes is shown next to each species acronym.
1718 Differentially enriched genes also correlated with PC1 Lifestyle (continuous proxy of
1719 industrialization) are shown in plain circles.

1720

1721 **Supplementary Figure 6 – Convergent elevation of gene-level Ka/Ks based on**
1722 **industrialization status.**

1723 Evidence for convergent positive selection on individual genes based on host industrialization
1724 status. The panel shows a heatmap of Ka/Ks differences for 90% gene families between
1725 industrialized and non-industrialized strains across species (x-axis). Genes with higher Ka/Ks
1726 values in industrialized strains are shown in purple, and those with higher values in non-
1727 industrialized strains are shown in green. Absent genes shown in white. Ka/Ks differences are
1728 displayed on a square-root-transformed color gradient. Mean Ka/Ks differences across species
1729 are shown on the left of the heatmap. Gene family annotations (preferred gene names or COG/KO
1730 IDs) are displayed on the y-axis. Genes with consistent Delta Ka/Ks signs across 8 or 7 of 10
1731 species are shown at the top and bottom, respectively.

1732

1733 **Supplementary Figure 7 – Convergent signals of SNV-host lifestyle associations across**
1734 **bacterial species at the level of KEGG KO categories.**

1735 Convergent signals of SNV-host lifestyle associations across bacterial species at the level of
1736 KEGG KOs. Each tile represents a KEGG KO that contains at least one 90% gene family with
1737 host industrialization-associated SNVs. KEGG KOs are shown along the y-axis, and species
1738 along the x-axis, cells are color-coded by species identity. All SNV data can be found in Supp.
1739 Table 7.