

1 **Human xenobiotic metabolism proteins have full-length and split**
2 **homologs in the gut microbiome**

3 Matthew Rendina^{1,2}, Peter J. Turnbaugh^{3,4}, Patrick H. Bradley^{1,2*}

4 ¹ Department of Microbiology, The Ohio State University, Columbus, OH 43210, USA

5 ² Infectious Diseases Institute, The Ohio State University, Columbus, OH 43210, USA

6 ³ Department of Microbiology and Immunology, University of California San Francisco, San
7 Francisco, California 94143, USA

8 ⁴ Chan-Zuckerberg Biohub-San Francisco, San Francisco, CA 94158, USA

9 * To whom correspondence should be addressed. Email: bradley.720@osu.edu

10 Running title: Gut microbial xenobiotic homologs

11 **Abstract**

12 Xenobiotics, including pharmaceutical drugs, can be metabolized by both host and microbiota,
13 in some cases by homologous enzymes. We conducted a systematic search for all human
14 proteins with gut microbial homologs. Because gene fusion and fission can obscure homology
15 detection, we built a pipeline to identify not only full-length homologs, but also cases where
16 microbial homologs were split across multiple adjacent genes in the same neighborhood or
17 operon (“split homologs”). We found that human proteins with full-length gut microbial
18 homologs disproportionately participate in xenobiotic metabolism. While this included many
19 different enzyme classes, short-chain and aldo-keto reductases were the most frequently
20 detected, especially in prevalent gut microbes, while cytochrome P450 homologs were largely
21 restricted to lower-prevalence facultative anaerobes. In contrast, human proteins with split
22 homologs tended to play roles in central metabolism, especially of nucleobase-containing
23 compounds. We identify twelve specific drugs that gut microbial split homologs may
24 metabolize; two of these, 6-mercaptopurine by xanthine dehydrogenase (XDH) and 5-
25 fluorouracil by dihydropyrimidine dehydrogenase (DPYD), have been recently confirmed in
26 mouse models. This work provides a comprehensive map of homology between the human
27 and gut microbial proteomes, indicates which human xenobiotic enzyme classes are most
28 likely to be shared by gut microorganisms, and finally demonstrates that split homology may
29 be an underappreciated explanation for microbial contributions to drug metabolism.

30 **Article Summary**

31 We develop a pipeline to systematically find human proteins with gut microbial homologs,
32 including those split across multiple microbial genes (e.g., operons). This reveals thousands of
33 proteins with full-length gut homologs, especially reductases and hydrolases that metabolize
34 xenobiotics. Nearly two dozen split homologs are also observed for central metabolic
35 enzymes, many of which can transform substrate analogs; in two cases, previous studies
36 verify that microbial split homologs enable the expected drug to be metabolized *in vivo*. These
37 results, which we provide as a resource, map out homology and shed light on parallel drug
38 metabolism between host and microbiome.

39 Introduction

40 Hundreds of small molecules, including drugs, can be metabolized by both human cells and
41 also the trillions of microorganisms that colonize the gastrointestinal tract (the gut microbiota)
42 [1,2]. In some cases (e.g. digoxin), drug metabolism by the microbiome can contribute to
43 observed differences in pharmacodynamics across patients. When drugs have a narrow
44 therapeutic window, microbial metabolism can be especially relevant, as even small
45 differences in concentration can lead to large changes in toxicity or efficacy [3]. However,
46 cases of microbial drug metabolism can be difficult to identify and are time- and labor-
47 intensive to characterize; for example, metabolism of digoxin by gut microbes was first
48 reported in 1981 [4], but the gene responsible was not identified until 2013 [5]. This has led to
49 interest in using bioinformatic methods, such as GutBug [6], MicrobeFDT [7], and SIMMER
50 [8], to help researchers prioritize the most promising gut microbial genes for further study.

51 Many proteins involved in drug metabolism are part of general “xenobiotic” systems, often with
52 relatively broad specificity, that transform or detoxify natural products. In humans, these
53 systems include cytochrome P450 proteins and glutathione-S-transferases. Many drugs are
54 derived from natural products that could be encountered in the environment. These natural
55 products typically exhibit high structural diversity, both because they are ammunition in “arms
56 races” between competitors, and because of other constraints on natural product enzyme
57 evolution [9]. We therefore might expect to find xenobiotic metabolic genes with less specific
58 substrate requirements in both hosts and microbiome. In contrast, other proteins involved in
59 drug metabolism have primary roles in central metabolism. These proteins typically have
60 narrower specificity and metabolize drugs that are structurally similar to their natural
61 substrates, regardless of whether they are found in nature, such as nucleoside analogs. Since
62 many central metabolic proteins are evolutionarily ancient, one might also expect to find cases
63 of direct homology between host and microbiome drug-metabolizing proteins. This may be
64 especially true for drugs like chemotherapeutic or immunomodulatory antimetabolites, as
65 these target conserved parts of metabolism.

66 Most approaches to detecting microbe-host homology focus on single genes. However,
67 horizontal transfer, multidomain protein architectures, and gene fusions can complicate this
68 picture [10–12], making it more difficult to determine whether orthologs of a human drug target

69 or drug-metabolizing enzyme are actually likely to exist in the microbiome. Eukaryotic
70 metabolic genes are especially likely to have bacterial origins [13], and many of these are
71 actually fusions of bacterial operons or domains (previously termed “S-genes”). Gene fusions
72 in eukaryotes may ensure co-expression in the absence of multi-cistronic operons, and may
73 also function to prevent metabolic intermediates from diffusing away in a larger cell volume
74 [12].

75 Taken together, this implies that gut microbial genomes may contain direct homologs to
76 human drug metabolism genes. Some of these may be part of more general systems, while
77 others may be central metabolic genes that happen to also metabolize designed substrate
78 analogs. While individual cases have been identified, the full extent of such homology has
79 remained unknown, as did whether particular systems or enzymes are particularly likely to be
80 shared across hosts and microbes. Furthermore, in some cases, these microbial homologs
81 may be encoded by multiple adjacent open reading frames; such “split” homologs would be
82 missed by a one-to-one homology search. Leveraging recently-published collections of
83 metagenome-assembled genomes (MAGs) from the human gut microbiome [14], we therefore
84 aimed to comprehensively identify gut microbial proteins that are either full-length or “split”
85 human homologs, then determine, based on curated human annotations, which of these were
86 most likely to participate in xenobiotic metabolism, and which specific roles those proteins are
87 most likely to play.

88 **Results**

89 **Thousands of human proteins have either full-length or split homologs in** 90 **the gut microbiome**

91 Our approach for identifying gut microbial homologs is described in **Figure 1**. Briefly, we
92 conducted a BLASTP homology search between gut microbial protein families and the human
93 proteome and identified cases where a human protein aligned to $\geq 2/3$ of a gut microbial
94 protein. To find full-length homologs, we kept the best microbial match per genome that also
95 aligned to $\geq 70\%$ of the human protein. To find split homologs, we identified sets of gut
96 microbial proteins that were jointly, but not individually, homologous to the majority of the

97 human protein and encoded by adjacent or near-adjacent genes on the same strand of the
98 same gut microbial assembly (see **Figure 1** and **Methods**).

99 Our gut microbial sequences came from the Universal Human Gut Proteome database
100 (UHGP) [14], which contains predicted protein sequences from >200K isolate and
101 metagenome-assembled genomes. Because UHGP contains a very large number of non-
102 identical protein sequences (>170M), we used a derivative of UHGP clustered at 90% amino
103 acid identity (UHGP-90), which retains most of the sequence diversity at less than one-tenth
104 the size (14M protein families) [15]. We then performed a local alignment search using
105 BLASTP [16] for each UHGP-90 protein cluster. Because we wanted to compare our results
106 against multiple databases, we did not limit our initial search to only known drug metabolism
107 genes, but instead searched against all >20K human protein sequences from UniProt [15],
108 and then filtered according to the criteria above.

109 Overall, we found that homology between the human and gut microbial proteomes was not
110 rare, with 12.9% of human proteins (2,589) having at least one gut microbial homolog.
111 Furthermore, while the majority had full-length homologs, a sizable minority (313) had at least
112 one split homolog. In fact, 23 human proteins had more split than full-length homologs, and 16
113 had no full-length homologs at all (**Figure 2, Table 1**), meaning that they could not have been
114 found by a conventional one-to-one homology search. These numbers are similar in
115 magnitude to a previous estimate of eukaryotic gene families that likely descended from
116 fusions of prokaryotic proteins or domains. This study identified 282 such families, 19 of which
117 were both widely distributed in eukaryotes and also “operon-like” in bacteria, in that they
118 appeared in an annotated operon in at least one bacterial genome [12].

119 **Human proteins with full-length vs. split homologs differ in function and** 120 **subcellular localization**

121 After identifying human proteins with full-length and split homologs, we used Gene Ontology
122 (GO) enrichment [17,18] to ask whether these two groups could be differentiated in
123 localization and function (**Supplementary Tables 1-2**). Considering proteins with mostly full-
124 length homologs, we observed that many of the enriched pathways were mitochondrial, such
125 as “tricarboxylic acid cycle” ($p_{adj} = 1.2 \times 10^{-13}$), “fatty acid beta-oxidation” ($p_{adj} =$

126 2.1×10^{-10}), and “carnitine metabolic process” ($p_{adj} = 1.9 \times 10^{-8}$). Because the eukaryotic
127 mitochondrion descends from a bacterial ancestor, we might expect human proteins with gut
128 bacterial homologs to localize to the mitochondrion. Indeed, human proteins with full-length
129 homologs were much more likely to localize to the mitochondrion (odds ratio 4.8, 95% CI [4.3,
130 5.5], $p < 2.2 \times 10^{-16}$, Fisher’s exact test). Further, this enrichment increased the more
131 frequently the full-length homologs were detected (**Supplementary Figure 1**). This set of
132 enrichments aligns strongly with previous work that identified a set of nuclear gene families
133 present in the last eukaryotic common ancestor that had mainly Alphaproteobacterial origins,
134 mitochondrial localization, and roles in energy production [19].

135 Remarkably, the most-enriched term among proteins with full-length homologs was
136 “xenobiotic metabolism” ($p_{adj} = 8.9 \times 10^{-25}$). There was an equally strong enrichment when
137 considering only non-mitochondrial genes. Further, the enrichment was not driven by a single
138 enzyme family. We observed homologs of short-chain and aldo-keto reductases,
139 carboxylesterases, arylamine N-acetyltransferases and arylacetamide deacetylases,
140 glutathione-S-transferases, flavin mono-oxygenases, UDP-glucuronosyl transferases and
141 cytochrome P450 family members, among others.

142 We next compared these results with the proteins that had mainly split homologs. In contrast,
143 these were not at all enriched for the “xenobiotic metabolism” GO term ($p_{adj} = 1$), but rather
144 for a smaller number of central pathways, namely purine, pyrimidine, and cofactor (folate and
145 molybdopterin) metabolism ($p_{adj} \leq 0.05$). Full-length homologs were also significantly
146 enriched for some of these pathways (**Supplementary Table 2**), but less so than the general
147 “xenobiotic metabolism” term, indicating that proteins with split homologs are a more
148 functionally specific group. Proteins in these central pathways, however, still make important
149 contributors to drug metabolism. Nucleoside and folate analogs, in particular, are common
150 antiviral, antibiotic, and chemotherapeutic agents.

151 Further, when we examined the subcellular distribution of human proteins with mainly split
152 homologs, the fraction localizing to the mitochondria was more modest, and did not differ
153 significantly from the base rate (odds ratio 2.0, 95% CI [0.23, 8.8], $p = 0.29$). If anything, the
154 proteins with the most split homologs were the least likely to be mitochondrial

155 **(Supplementary Figure 1)**. Proteins with split homologs therefore appear to participate in
156 different biological processes (cytosolic, primarily central metabolism) than full-length
157 homologs (mitochondrial, energy production, both xenobiotic and central metabolism).

158 **Reductases and hydrolases dramatically outnumber cytochromes and** 159 **UDP-glucuronosyltransferases in gut microbes**

160 The above analysis indicates the presence of full-length homologs of xenobiotic metabolism
161 enzymes in gut microbes. However, it does not tell us about their phylogenetic distribution,
162 which is important because gut microbial clades vary in their prevalence and average
163 abundance across orders of magnitude [20]. We therefore identified eleven enzyme families
164 with at least some members known to participate in human xenobiotic metabolism, then
165 determined which gut microbial species contained homologs of these families (**Figure 2,**
166 **Table 2**), as well as how many distinct bacterial proteins were identified (**Table 3**).

167 While cytochrome P450s are one of the most important and well-studied xenobiotic
168 detoxification systems in humans, we saw relatively few homologs in gut microbes. Their
169 homologs were also mainly restricted to facultative anaerobes, which makes sense given the
170 oxygen-dependent mechanisms of these proteins. Furthermore, the species with the most
171 cytochrome P450 homologs were in low-prevalence families like the *Paenibacillaceae* (**Table**
172 **2**). This suggests that while cytochrome P450 homologs can be found, they may be less likely
173 to be relevant to the adult human gut, though with the caveat that gut oxygenation can also
174 vary across development and in disease [21]. Flavin-dependent monooxygenases and UDP-
175 glucuronosyltransferases had similarly sparse distributions in mostly lower-abundance gut
176 microbes.

177 Two enzyme types had intermediate distributions. First, glutathione S-transferase homologs
178 were much more commonly detected than cytochrome P450s, but were almost exclusively
179 found in Proteobacterial facultative anaerobes, like the *Enterobacteriaceae* and
180 *Burkholderiaceae*. However, while these are typically low-abundance, *Enterobacteriaceae* are
181 prevalent, and Proteobacteria can rise to high levels in certain individuals and situations,
182 making them major contributors to functional variability [22]. This suggests that gut microbial
183 GST activity might also be especially variable across individuals. Second, arylamine

184 acetylases were most observed in facultative anaerobes (*Enterobacteriaceae* and
185 *Staphylococcaceae*), yet were also detected in certain *Lachnospiraceae*, the most prevalent
186 gut microbial family worldwide [20]. Substrates for these genes include the anti-hypertensive
187 vasodilator hydralazine [23] and the anti-tubercular isoniazid [24]. Interestingly, it has been
188 previously shown that isoniazid is also metabolized by strains of *M. tuberculosis* by an
189 arylamine acetylase homologous to human NAT2 [25]. Finally, arylamine acetylases are also
190 responsible for both increasing and decreasing the carcinogenicity of certain environmental
191 pollutants, suggesting that gut microbes could also modulate these risks [26].

192 In contrast, we detected thousands of short-chain and aldo-keto reductases in common gut
193 microbes, like *Lachnospiraceae*, *Enterobacteriaceae*, and *Bacteroidaceae* (**Table 3**). In
194 humans, both classes of enzymes act on a wide range of substrates; notably, certain
195 members can participate in the reduction of steroid-like and polycyclic molecules, including
196 bile acid intermediates [27,28]. Gut microbes are known for their ability to transform primary
197 bile acids into secondary bile acids, and this metabolism has well-studied consequences for
198 immune and metabolic signaling in the host [29]. Additionally, the *Lachnospiraceae* member
199 *Clostridium bolteae* was recently found to directly metabolize the steroids nabumetone,
200 hydrocortisone, and tacrolimus via the gene DesE [30]. We detected that the UHGP-90
201 protein with the best hit to DesE (GUT_GENOME228173_01934) appeared to be a full-length
202 homolog of the human protein PEGR, a trans-2-enoyl-CoA reductase that is a member of the
203 SDR family. The prominence of reductases in the most common gut microbes aligns with
204 previous observations that reduction reactions are especially common ways for gut microbes
205 to transform xenobiotics, potentially because of the need for alternative electron acceptors in
206 the absence of molecular oxygen [31,32].

207 Homologs of two other redox-active enzyme classes, aldehyde dehydrogenases and quinone
208 oxidoreductases, were also observed frequently in *Lachnospiraceae*, but the highest number
209 of distinct bacterial homologs were found in facultative anaerobes like *Enterobacteriaceae*,
210 *Lactobacillaceae*, or *Burkholderiaceae*. Enzymes in these families, of course, play roles in
211 both central and xenobiotic metabolism, complicating their interpretation. For example,
212 aldehyde dehydrogenase oxidizes acetaldehyde to acetate (or the reverse, in microbial
213 ethanol production). However, aldehyde dehydrogenase enzymes can have a variety of other

214 substrates (e.g. lactaldehyde [33]) and are also involved in the detoxification of drugs like
215 cyclophosphamide [34].

216 Finally, homologs of type B carboxylesterases and the “GDYG” group of lipases (which
217 include hormone-sensitive lipases, arylacetamide deacetylases, and neutral cholesterol
218 esterases) were also found frequently, especially in *Lachnospiraceae* and *Bacteroidaceae*. In
219 humans, in addition to deactivating drugs like flutamide [35] and indiplon [36], these
220 hydrolases bioactivate a large number of prodrugs, including enalapril [37] and irinotecan [38].
221 Overall, this analysis shows that while several systems used by humans to detoxify
222 pharmaceutical, dietary, and environmental compounds do have at least some analog in the
223 gut microbiome, certain enzyme families are much better represented in the most prevalent
224 gut microbes. Specifically, these include redox-active enzymes, especially short-chain and
225 aldo-keto reductases, and hydrolases, including lipases and carboxylesterases.

226 **Identifying split and full-length homologs of specific drug-metabolizing** 227 **genes**

228 In many cases, we know the specific substrates on which drug-metabolizing enzymes act. We
229 therefore used the database PharmGKB [39] to identify cases where human proteins with gut
230 homologs were known to be involved in the metabolism of either a pharmaceutical drug or one
231 of its downstream metabolites.

232 Out of 154 proteins in PharmGKB with reviewed entries in UniProt, we found that a large
233 majority (126/154, 82%) had at least one full-length or split homolog. 97% of these (122/126)
234 had more full-length than split gut homologs; this set of proteins metabolized 215 drugs in total
235 (**Supplementary Table 3**). Consistent with the sparse distribution we observed above,
236 cytochromes were the most under-represented category in this list, with only eight genes
237 found to have gut microbial homologs compared to 22 in PharmGKB as a whole. In contrast,
238 12 out of 13 drug-metabolizing UDP-glucuronosyl-transferases and all nine aldo-keto
239 reductases were found to have full-length gut homologs.

240 Interestingly, despite the strong enrichment we observed for xenobiotic metabolism among
241 proteins with full-length homologs, <50% of these proteins (56/122) fell into one of the ten
242 classes listed above. While many different types of enzymes were represented among the

243 remainder, a plurality of 36 were annotated in GO as metabolizing nucleobase-containing
244 compounds. This is consistent with the observation that this process was enriched among
245 both full-length and split homologs. Furthermore, nucleobase-containing analogs are some of
246 the most common human chemotherapeutic, immunomodulatory, and antiviral drugs, and
247 their metabolism is also well-studied, as variants that affect their metabolism have large
248 consequences for health. Finally, of the remaining 30 genes, 20 were annotated in GO as
249 oxidoreductases, further underscoring the importance of redox-active genes in gut microbial
250 metabolism.

251 When we instead kept only cases with more split than full-length homologs, we found four
252 genes involved in the metabolism of 12 drugs (**Table 4**). Again, three of these genes were
253 involved in nucleotide metabolism, and many of these drugs were antimetabolite
254 chemotherapeutics such as thioguanine, doxorubicin, and mercaptopurine. We noted that two
255 out of four genes metabolized 5-fluorouracil (dihydropyrimidine dehydrogenase, or DPYD; and
256 uridine monophosphate synthase, or UMPS). These results indicate that both full-length and
257 split gut homologs may play roles in the microbial transformation of nucleoside analogs.

258 The genes dihydropyrimidine dehydrogenase (DPYD) and xanthine dehydrogenase (XDH)
259 had among the most split homologs. In the case of XDH, human gut microbes have been
260 shown to catabolize purines such as uric acid, an endogenous substrate of human XDH, and
261 to alter purine levels *in vivo* in a mouse model; knockout experiments suggest that both
262 phenomena require an operon bearing an XDH homolog [40]. This potential conservation of
263 function supports a potential role for XDH homologs in the metabolism of purine analog drugs,
264 such as azathioprine (AZA), a chemotherapeutic and immunosuppressive drug that is given
265 orally. Indeed, it has recently been shown in an *in vivo* preclinical model that *Blautia wexlerae*
266 reduces the therapeutic effect of AZA by metabolizing its active metabolite, 6-mercaptopurine
267 (6MP) into the inactive form, 6-TX; furthermore, this metabolism can be interrupted with the
268 XDH inhibitor allopurinol [41].

269 Recent publications also support that both the human DPYD protein and its bacterial
270 counterparts, PreT and PreA (encoded by the *preTA* operon) can inactivate the
271 chemotherapeutic drug 5-fluorouracil (5-FU) *in vivo*. While 5-FU itself is not given orally, the
272 orally-available prodrug, capecitabine, can also be activated to 5-FU by host liver enzymes as

273 well as select gut bacterial strains [2]. Indeed, in a mouse model of colorectal cancer
274 treatment with 5-FU, mice monocolonized with a *preTA* knockout strain of *E. coli* had better
275 survival than those monocolonized with a *preTA* overexpression strain [2].

276 **Discussion**

277 We conducted a systematic survey of homology between the human and gut microbial
278 proteomes. This analysis included both full-length and “split” homologs. We found that around
279 one in ten human proteins (2.6K) had at least some homolog in the gut microbial proteome,
280 and that 23 human proteins had primarily split homologs. While our focus was on drug and
281 xenobiotic metabolism, such a map of host-microbial homology may also be helpful to
282 microbiome researchers more broadly. With this in mind, our code and results are available as
283 publicly available resources (see **Data Availability**).

284 Among human proteins with full-length homologs, xenobiotic metabolism was the most
285 enriched process, and many different xenobiotic enzyme classes were found in gut genomes.
286 However, the most predominant systems in humans (cytochrome P450s, glutathione-S-
287 transferases) were relatively rare in the gut microbiome. Instead, reductases and hydrolases
288 were the most common, especially among the most prevalent microorganisms. This is
289 consistent with previous observations about types of drug metabolism engaged in by the gut
290 microbiome [31,32], and builds on these observations by enumerating specific classes of
291 enzymes that are likely to contribute. Proteins involved in central metabolism, especially of
292 nucleoside-containing compounds, were also commonly found as both full-length and split
293 homologs in the gut microbiome. Finally, in two cases (6MP and 5-FU), the gut microbial split
294 homologs we identified have been shown to metabolize pharmaceutical nucleoside analogs in
295 mouse models [2,41].

296 One limitation of this work is that we have only considered the gut microbiome. This
297 community was our focus because microbial biomass is highest in the gut [42], because orally
298 ingested drugs are absorbed in the intestine [43], and because the gut is closely connected to
299 the liver, the primary site of drug metabolism in humans [44]. However, split orthologs in skin

300 microbes may also be relevant for topically applied drugs, and similarly for oral microbes and
301 drugs delivered as rinses.

302 A technical limitation of this work is that sequencing and annotation errors can give rise to *in*
303 *silico*, artifactual gene “fusions” or “fissions.” We believe that the way that UHGP-90 protein
304 clusters were constructed would favor such “fusions.” UHGP-90 protein clusters were
305 constructed using MMSeqs2’s “includst” algorithm [45] in target-coverage mode, meaning that
306 the representative sequence for a cluster must cover 80% of each member sequence, but not
307 necessarily vice versa. This has the advantage that protein fragments or artifactual “fissions”
308 would seldom be chosen as representative sequences, but also means artifactual “fusions”
309 would be chosen more often. Since we use the representative sequences in this pipeline, this
310 effect would therefore bias us away from detecting split homologs.

311 Of course, it is important to emphasize that homologs may differ in substrate specificity. This
312 is especially true over long evolutionary distances (e.g., between humans and
313 microorganisms) and for enzymes whose substrate specificity is broad (e.g. many xenobiotic
314 metabolism genes). Follow-up experiments would therefore be necessary to establish whether
315 specific substrates are shared between host and microbial homologs. Advances in
316 computational structural biology, such as improvements to high-throughput ligand docking
317 tools [46–48], may also help prioritize homologs that could contribute to parallel drug
318 metabolism between host and microbiome. We speculate that interactions between central
319 metabolic enzymes and substrate analogs, such as the chemotherapeutics 6MP and 5-FU,
320 may be especially likely to translate: these enzymes are more evolutionarily constrained than
321 broad-spectrum xenobiotic enzymes [9], and the corresponding drugs bind in the active site,
322 which is typically highly conserved.

323 While we have focused on proteins that metabolize drugs, the protein targets of drugs could
324 also be conserved, potentially causing off-target effects on the microbiome. Such unintended
325 effects of pharmaceuticals on gut microorganisms are not rare: one study of more than 1,000
326 marketed drugs found that nearly a quarter inhibited the growth of at least one of 40
327 representative gut isolates [49]. As above, we would expect host-microbiome homology to be
328 especially relevant when considering proteins targeted by substrate analogs. Indeed, a study

329 of the chemotherapeutic 5-FU showed that it had large effects on gut microbial growth [2], and
330 antimetabolites as a class were also enriched for antimicrobial effects in the study above [49].

331 A final unresolved question is the evolutionary history of these homologs. For example, which
332 of these homologs are descendants of ancestral sequences present in the last universal
333 ancestor, and which might be better explained by horizontal gene transfer? Transfers from
334 bacteria into early eukaryotes [12], especially from the ancestors of modern organelles [50],
335 as well as transfers from modern eukaryotes into bacteria [51], are two mechanisms that could
336 lead to both full-length and split homology. While ancient events are intrinsically difficult to
337 resolve, phylogenetic methods, combined with the current explosion in microbial genome
338 sequencing, may enable us to distinguish between these possibilities.

339 **Methods**

340 **Identification of split and full-length homologs**

341 To identify homologs, we performed a BLASTP search of all 13.9M proteins in the UHGP-90
342 database against all 20.6K human proteins downloaded from UniProt (2023/9/13). This
343 yielded 8.5M potential matches. We then performed the following filtering steps:

- 344 1. Best human hit: for each UHGP-90 protein, retain only the human protein with the
345 highest bitscore;
- 346 2. Microbe coverage: retain only alignments covering at least 67% of the prokaryotic
347 UHGP-90 protein sequence; additionally, retain only UHGP-90 sequences that are at
348 least 80 amino acids long.

349 The pipeline diverged after this point for full-length and split homologs. For full-length
350 homologs, we were interested in individual microbial proteins where an alignment covered
351 most of the human protein, and where this alignment was unlikely to be due to contamination
352 in the microbial genomes. We therefore performed the following filtering steps:

- 353 F3. Human coverage: each alignment must cover at least 70% of the human sequence;

354 F4. Contamination: filter out any alignments whose amino acid percent ID was more than
355 three standard deviations above the mean (mean: 30% ID; cutoff: 51.8% ID).

356 For split homologs, we sought to determine which UHGP-90 families were encoded by
357 neighboring features in the same genome, and where this was unlikely to be the result of
358 contamination, as above. Because multiple genomes could encode the same UHGP-90
359 family, we had to first expand our results, then filter, as follows:

360 S3. Same genome: first, determine which individual UHGG genomes encoded multiple
361 UHGP-90 families aligning to each individual human protein. Then, retain only those
362 alignments, repeated for each genome that encoded the UHGP-90 protein;

363 S4. Joint human coverage: for each UHGG genome and for each human protein,
364 determine whether the alignments between the human protein and the UHGP-90
365 proteins from that genome could jointly, but not individually, cover at least 70% of the
366 human sequence;

367 S5. Feature distance: for each UHGG genome and for each human protein, compute the
368 minimum distance (in feature numbers) between each UHGP-90 protein aligning to
369 that human protein. Retain only sets of alignments with at least two different UHGP-
370 90 proteins that are three or fewer features apart, and that are additionally all on the
371 same strand and contig. Then, repeat the human coverage step to ensure that the
372 remaining alignments still jointly cover $\geq 70\%$ of the human protein, as some have
373 been removed;

374 S6. Contamination: filter out any alignments whose amino acid percent ID was more than
375 three standard deviations above the mean observed for all full-length homologs, as
376 above (mean: 30% ID; cutoff: 51.8% ID); repeat the human coverage step again and
377 report results.

378 Filtering and analysis steps were carried out in a Snakemake pipeline, using Pandas [52],
379 Polars [53], and R with Tidyverse [54,55].

380 **Enrichment analysis**

381 Gene Ontology (GO) annotations [17] from UniProt [15] were used to determine subcellular
382 localization (“cellular component”) and function (“biological process”) for human proteins.
383 Proteins whose cellular component annotations matched the regular expression
384 “[Mm]mitochondr” were retained as mitochondrially-localized. Enrichment analysis was carried
385 out using TopGO [56] using Fisher’s exact test on GO biological process terms, with the
386 resulting *p*-values corrected for multiple testing using the Benjamini-Hochberg method [57].

387 **Analyzing subclasses of xenobiotic enzyme families**

388 Xenobiotic enzyme families were defined using protein family annotations in UniProt, using
389 regular expression matches for the following:

- 390 • Aldo-keto reductases (“akr”): “Aldo/keto reductase family”;
- 391 • UDP-glucuronosyltransferases (“udp”): “UDP-glycosyltransferase family” (note: all human
392 members of this family except cerebroside synthase were annotated as UDP-
393 glucuronosyltransferases);
- 394 • Glutathione S-transferases (“gst”): “GST superfamily” (note: this included the alpha, zeta,
395 sigma, pi, mu, theta, omega, and kappa families);
- 396 • Arylamine N-acetyltransferases (“aryl”): “Arylamine N-acetyltransferase family”;
- 397 • GDXG-like hydrolases (“gdxg”): “GDXG’ lipolytic enzyme family”;
- 398 • Cytochrome P450s (“cyto”): “Cytochrome P450 family”;
- 399 • Type B carboxylesterases (“ester”): “Type-B carboxylesterase/lipase family”;
- 400 • Flavin monooxygenases (“flavin”): “Flavin monoamine oxidase family|FMO family” (note:
401 this included the FIG1 subfamily);
- 402 • Short-chain reductases (“sdr”): “Short-chain dehydrogenases/reductases (SDR)”;
- 403 • Quinone oxidoreductases (“quin”): “Quinone oxidoreductase subfamily”.

404 Full-length homologs were partitioned into one of the above classes. Next, for each class, the
405 phylogenetic diversity of species containing at least one full-length homolog was calculated
406 using Faith's PD [58]. The phylogeny used was the maximum-likelihood tree of the 4,616
407 species in UHGG [14] generated via IQ-TREE [59], which we midpoint-rooted using APE [60].
408 Faith's PD was calculated using Picante [61]. Xenobiotic classes were visualized in
409 descending order of PD. The number of unique species with at least one full-length homolog
410 is given in **Table 2**, while the total number of unique UHGP-90 IDs per family is given in **Table**
411 **3**. Results were visualized using the R package ggtree [62].

412 **Identification of gut homologs of drug-metabolizing enzymes**

413 Pathway, gene, relationship, and chemical annotations were downloaded from PharmGKB
414 (2024/10/12) [63]. HUGO Gene Nomenclature Committee (HGNC) identifiers [64] in
415 PharmGKB were mapped, using data downloaded from HGNC (2024/10/01), to UniProt IDs.
416 Chemicals of interest in PharmGKB were defined as having the chemical classes "Drug",
417 "Drug Class", "Prodrug", or "Metabolite" (this refers to drug metabolites, not endogenous
418 substrates or "Biological Intermediates", which were excluded). This was necessary because
419 certain endogenous human metabolites were annotated in pathways, but only peripherally
420 related to drug metabolism, e.g., homocysteine in the methotrexate metabolism pathway
421 (present because methotrexate targets folate biosynthesis, which in turn is linked to
422 homocysteine via the S-adenosyl-methionine cycle). Reactions from all PharmGKB pathways
423 were then filtered such that the reactant and products were different, the reaction type was not
424 "Transport", the "Controller" (typically an enzyme or regulator) was known, and either the
425 reactant or product (or both) was a chemical of interest as defined above.

426 **Identification of other xenobiotic enzyme types in PharmGKB**

427 In addition to the xenobiotic enzyme classes we defined above, we also considered two other
428 classes of enzymes:

- 429 • Nucleobase metabolism genes: genes annotated to the GO term "nucleobase-containing
430 compound metabolic process (GO: 0006139)", or any term below it, but excluding genes
431 annotated to the following GO terms or any terms below them:

- 432 ○ “Acyl-CoA metabolic process (GO:0006637)”;
- 433 ○ “Coenzyme A metabolic process (GO:0015936)”;
- 434 ○ “FMN metabolic process (GO:0046444)”;
- 435 ○ “FAD metabolic process (GO:0046443)”;
- 436 ○ “Pyridine nucleotide metabolic process (GO:0019362)”.
- 437 ● Oxidoreductases: genes annotated to the GO term “oxidoreductase activity
- 438 (GO:0016491)”.

439 **Identification of DesE homologs**

440 In the study showing DesE from *Clostridium bolteae* metabolized the ketone group of
441 pharmaceutical steroids [30], its GenBank [65] protein accession was given as EDP16280.1.
442 To determine whether this protein was represented in our full-length homologs, we first
443 retrieved this accession and used it to perform an MMSeqs2 [45] search in “easy-search”
444 mode against all UHGP-90 protein sequences. The best-hit protein
445 (GUT_GENOME228173_01934, 97.7% identity) was then used to filter our list of full-length
446 homologs. GUT_GENOME228173_01934 was found to be a full-length homolog of the human
447 gene PECR (UniProt ID Q9BY49), and was distributed in six Lachnospiraceae species; these
448 included *Clostridium_M bolteae* where it was discovered and two other species in the genus
449 *Clostridium_M*.

450 **Data availability**

451 The UHGG and UHGP datasets [14] can be obtained from EMBL-EBI [66] at
452 https://ftp.ebi.ac.uk/pub/databases/metagenomics/mgnify_genomes/human-gut/v1.0/. The
453 human reference proteome was downloaded from UniProt [15] and is available at
454 https://ftp.uniprot.org/pub/databases/uniprot/previous_releases/release-2023_03/. HGNC IDs
455 [64] can be obtained at [https://storage.googleapis.com/public-download-](https://storage.googleapis.com/public-download-files/hgnc/archive/archive/monthly/tsv/hgnc_complete_set_2024-10-01.txt)
456 [files/hgnc/archive/archive/monthly/tsv/hgnc_complete_set_2024-10-01.txt](https://storage.googleapis.com/public-download-files/hgnc/archive/archive/monthly/tsv/hgnc_complete_set_2024-10-01.txt).

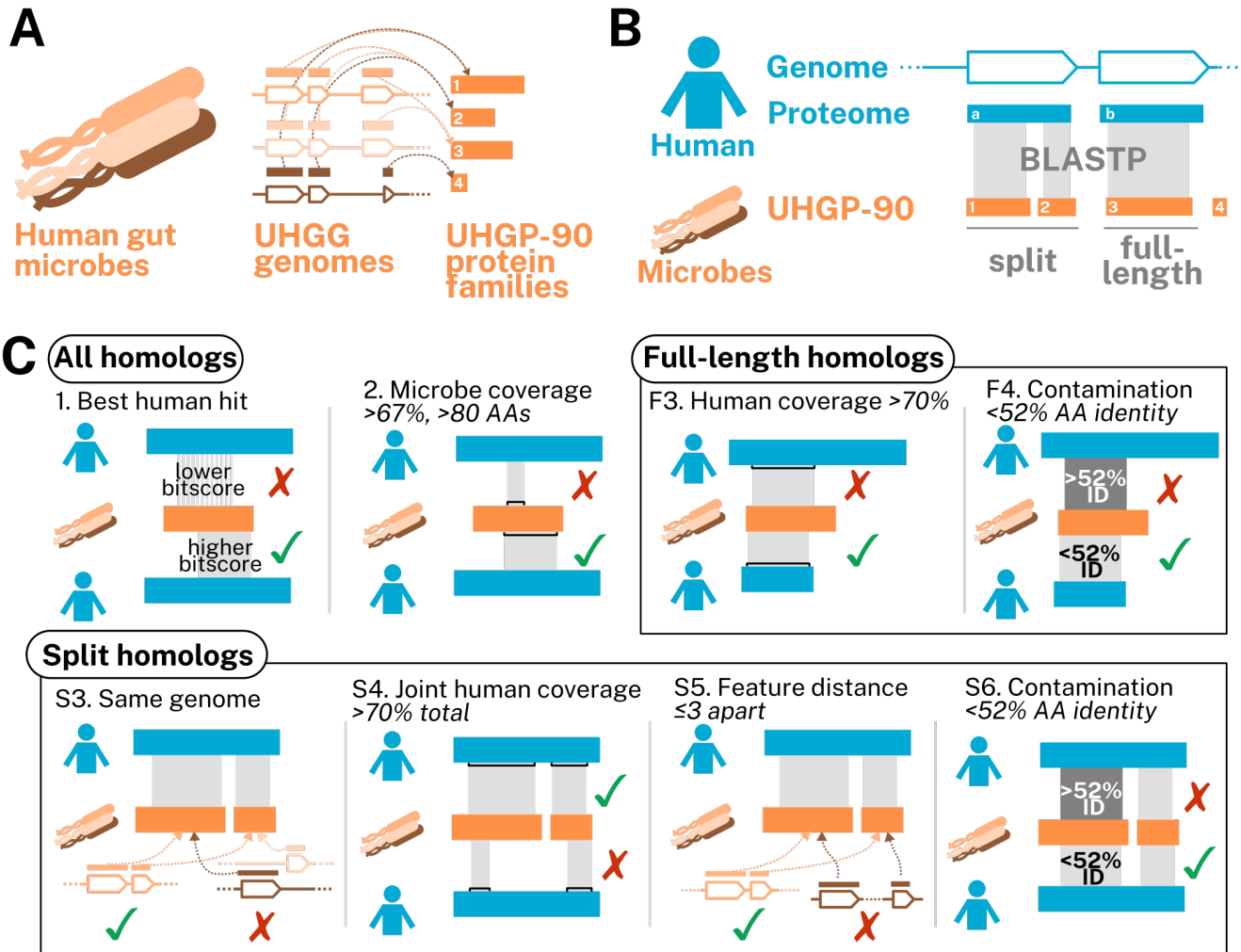
457 Code for our pipeline and analysis are available at
458 https://github.com/pbradleylab/split_homology. The processed output of the pipeline
459 (described starting on p. 13) is available via Zenodo at <https://zenodo.org/uploads/14037045>
460 (DOI: 10.5281/zenodo.14037045).

461 **Acknowledgements**

462 The authors wish to thank Abigail Lind for helpful feedback. Funding was provided by The
463 Ohio State University (P.H.B.) and the National Institutes of Health: R35GM151155 (P.H.B.);
464 R01CA255116, R01HL122593 (P.J.T). P.J.T is a Chan Zuckerberg Biohub-San Francisco
465 Investigator.

466

Figures and Tables



467

468

469

470

471

472

473

474

475

476

477

478

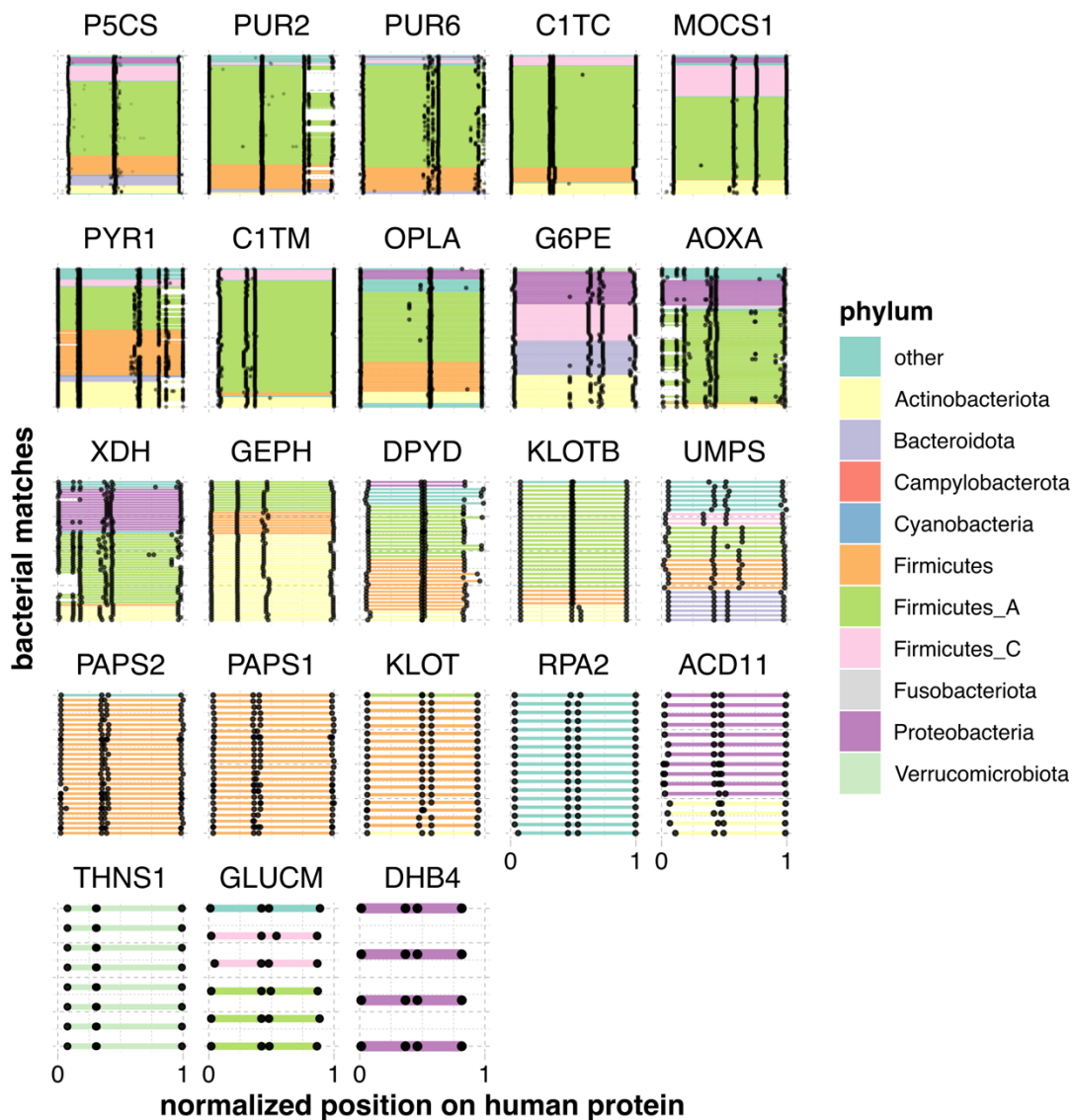
479

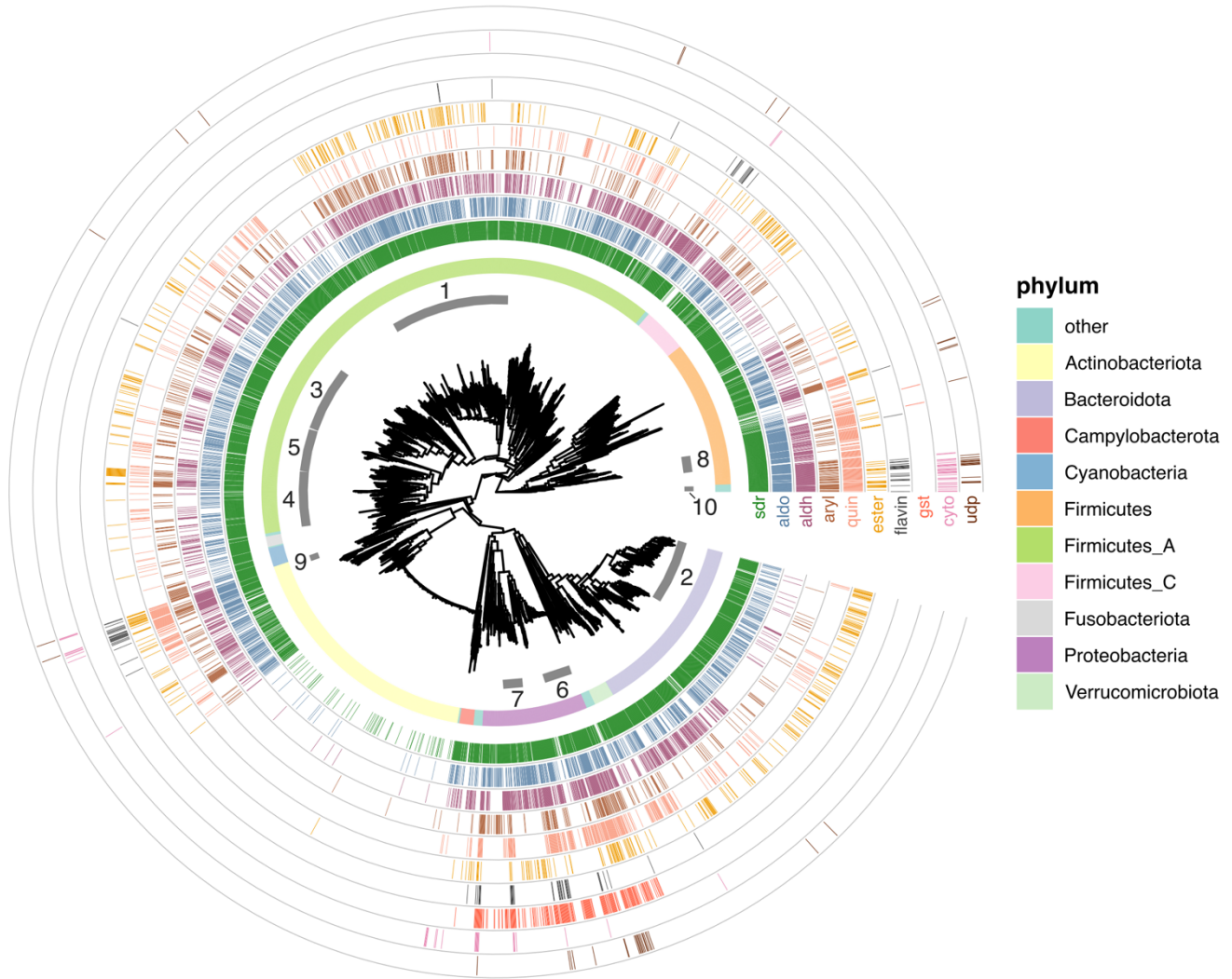
480

481

482

Figure 1: Schematic showing an overview of our approach. A) Diagram showing microbial gene and protein catalog. Human gut microbial genomes are collected in UHGG. Open reading frames (outlines, multiple shades) encode protein sequences (solid orange blocks, multiple shades). These are clustered at 90% amino acid identity (dashed lines) to form the UHGP-90 protein families (solid orange blocks, single shade). B) Pipeline overview. BLASTP is used to identify cases where microbial UHGP-90 proteins jointly or individually align to human proteins. Proteins that jointly align and are encoded by nearby features on the same genome are termed “split homologs.” C) Filtering criteria. Steps 1 and 2 involve finding human proteins (“best human hit”) that align to at least two-thirds of a microbial protein (“microbe coverage”). These steps are common to both pipelines (“all homologs”). For full-length homologs, alignments are then filtered based on whether they individually cover $\geq 70\%$ of a human protein (“human coverage”) at less than 52% amino acid identity (“contamination”). For split homologs, bacterial proteins aligning to a human protein must be encoded in the same genome (“same genome”), cover $\geq 70\%$ of the human protein (“joint human coverage”), be at most three features apart (“feature distance”), and have <52% amino acid identity (“contamination”). Note that in the “split homologs” section, step 3 is re-run after steps 4, 5, and 6. See **Methods** for details.





492

493 **Figure 3: Phylogenetic distribution of full-length gut microbial xenobiotic homologs.** The inset is a
494 midpoint-rooted tree of all bacteria included in UHGG v1.0. Numbered ring segments indicate selected
495 bacterial families (see Tables 2-3). The first complete ring shows the phylum-level classification. Note that
496 “Firmicutes” contains Bacilli, “Firmicutes_A” contains Clostridia, and “Firmicutes_C” contains Negativicutes.
497 Successive ring tracks mark the species where full-length homologs were identified (colored lines). From
498 inner to outer, these are: short-chain reductases (“sdr”, dark green); aldo-keto reductases (“aldo”, blue-
499 green); aldehyde dehydrogenases (“aldh”, purple); arylamine and arylacetamide metabolism (“aryl”, brown);
500 quinone oxidoreductases (“quin”, salmon); type B carboxylesterases (“ester”, orange), flavin-containing
501 mono-oxygenases (“flavin”, dark gray), glutathione-S-transferases (“gst”, red), cytochromes (“cyto”, violet),
502 and UDP-glucuronosyl-transferases (“udp”, dark brown). Ring tracks are plotted in order of Faith’s
503 phylogenetic diversity (PD), descending from inside to outside. Numbered families: 1. *Lachnospiraceae*,
504 2. *Bacteroidaceae*, 3. *Oscillospiraceae*, 4. *Acutalibacteriaceae*, 5. *Ruminococcaceae*,
505 6. *Enterobacteriaceae*, 7. *Burkholderiaceae*, 8. *Lactobacillaceae*, 9. *Mycobacteriaceae*, 10.
506 *Paenibacillaceae*.

507 **Tables**

Entry	Entry Name	Description	nFull	nSplit
P54886	P5CS_HUMAN	Delta-1-pyrroline-5-carboxylate synthase	27	293
P22102	PUR2_HUMAN	Trifunctional purine biosynthetic protein adenosine-3	0	200
P22234	PUR6_HUMAN	Bifunctional phosphoribosylaminoimidazole carboxylase/phosphoribosylaminoimidazole succinocarboxamide synthetase	0	145
P27708	PYR1_HUMAN	Multifunctional protein CAD	0	120
Q9NZB8	MOCS1_HUMAN	Molybdenum cofactor biosynthesis protein 1	0	75
P11586	C1TC_HUMAN	C-1-tetrahydrofolate synthase, cytoplasmic	0	64
O14841	OPLA_HUMAN	5-oxoprolinase	7	54
O95479	G6PE_HUMAN	GDH/6PGL endoplasmic bifunctional protein	7	47
Q6UB35	C1TM_HUMAN	Monofunctional C1-tetrahydrofolate synthase, mitochondrial	4	47
Q06278	AOXA_HUMAN	Aldehyde oxidase	3	46
P47989	XDH_HUMAN	Xanthine dehydrogenase/oxidase	0	38
O95340	PAPS2_HUMAN	Bifunctional 3'-phosphoadenosine 5'-phosphosulfate synthase 2	0	28
Q9NQX3	GEPH_HUMAN	Gephyrin	0	27
Q12882	DPYD_HUMAN	Dihydropyrimidine dehydrogenase	24	26
O43252	PAPS1_HUMAN	Bifunctional 3'-phosphoadenosine 5'-phosphosulfate synthase 1	0	24
P11172	UMPS_HUMAN	Uridine 5'-monophosphate synthase	0	22
Q86Z14	KLOTB_HUMAN	Beta-klotho	0	12
Q709F0	ACD11_HUMAN	Acyl-CoA dehydrogenase family member 11	0	11
Q7Z3D6	GLUCM_HUMAN	D-glutamate cyclase, mitochondrial	0	6
Q9H9Y6	RPA2_HUMAN	DNA-directed RNA polymerase I subunit RPA2	2	5
Q9UEF7	KLOT_HUMAN	Klotho	0	5
P51659	DHB4_HUMAN	Peroxisomal multifunctional enzyme type 2	0	1
Q8IYQ7	THNS1_HUMAN	Threonine synthase-like 1	0	1

508 **Table 1: Table showing all human proteins with more split than full-length homologs (nSplit, nFull)**
 509 **in gut bacteria.**

510

511

phylum	family	sdr	aldo	aldh	gdxg	quin	ester	aryl	flavin	gst	cyto	udp
Firmicutes_A	Lachnospiraceae	192	178	170	129	34	105	27	4	0	1	0
Bacteroidota	Bacteroidaceae	116	91	26	18	26	67	4	0	0	0	0
Proteobacteria	Enterobacteriaceae	108	105	108	63	107	55	62	8	107	1	19
Firmicutes	Lactobacillaceae	59	59	49	11	48	3	0	0	0	0	0
Firmicutes_A	Ruminococcaceae	54	47	34	26	10	21	5	0	0	0	0
Firmicutes_A	Acutalibacteraceae	45	35	31	21	18	14	2	0	0	0	0
Firmicutes_A	Oscillospiraceae	43	36	37	18	11	17	11	1	1	0	1
Firmicutes_A	Clostridiaceae	31	22	30	12	5	2	0	8	0	3	1
Campylobacterota	Campylobacteraceae	25	6	21	1	0	1	0	0	0	6	0
Proteobacteria	Burkholderiaceae	22	16	22	9	13	12	2	7	21	4	0
Firmicutes	Enterococcaceae	20	20	18	12	20	5	3	1	1	0	0
Actinobacteriota	Mycobacteriaceae	19	20	20	10	20	17	2	9	0	2	0
Firmicutes_I	Paenibacillaceae	19	19	18	12	19	12	6	2	0	7	3
Firmicutes	Staphylococcaceae	14	14	14	3	13	10	14	1	0	0	3
Proteobacteria	Moraxellaceae	11	9	11	9	11	0	1	8	11	0	0
Bacteroidota	Tannerellaceae	11	11	10	2	7	5	2	0	0	0	4
Proteobacteria	Pseudomonadaceae	10	9	10	6	10	1	2	7	10	1	3
Firmicutes	Amphibacillaceae	7	7	7	0	7	1	3	1	0	6	0
Firmicutes	Bacillaceae	4	4	4	0	4	4	2	4	0	4	4
Firmicutes	Bacillaceae_G	4	4	4	3	4	0	4	4	0	3	4
Proteobacteria	Xanthobacteraceae	3	3	3	3	3	1	2	3	3	3	1

512 **Table 2: Number of unique species per family in UHGG (rows) where at least one homolog in a**
 513 **particular xenobiotic enzyme class (columns) was detected.** Highlighted cells show the top three
 514 families per xenobiotic class.

515

516

phylum	family	sdr	aldo	aldh	gdxg	quin	ester	aryl	flavin	gst	cyto	udp
Firmicutes_A	Lachnospiraceae	3,853	910	583	231	54	253	33	4	0	1	0
Bacteroidota	Bacteroidaceae	2,037	498	66	23	28	171	5	0	0	0	0
Proteobacteria	Enterobacteriaceae	1,584	177	669	57	278	42	40	7	186	1	15
Firmicutes	Lactobacillaceae	588	309	76	12	194	3	0	0	0	0	0
Firmicutes_A	Ruminococcaceae	1,403	246	108	65	20	87	9	0	0	0	0
Firmicutes_A	Acetivibacteraceae	1,162	286	124	68	47	84	8	0	0	0	0
Firmicutes_A	Oscillospiraceae	1,653	259	171	48	26	52	18	1	1	0	1
Firmicutes_A	Clostridiaceae	338	60	91	12	6	2	0	12	0	2	1
Campylobacterota	Campylobacteraceae	86	16	12	1	0	1	0	0	0	4	0
Proteobacteria	Burkholderiaceae	729	191	238	20	60	43	3	7	123	5	0
Firmicutes	Enterococcaceae	222	75	34	12	83	4	2	1	1	0	0
Actinobacteriota	Mycobacteriaceae	267	35	127	19	54	20	2	13	0	23	0
Firmicutes_I	Paenibacillaceae	595	63	64	18	65	15	6	2	0	11	3
Firmicutes	Staphylococcaceae	168	39	68	3	26	9	14	2	0	0	2
Proteobacteria	Moraxellaceae	201	8	104	18	32	0	1	23	27	0	0
Bacteroidota	Tannerellaceae	542	121	34	3	8	23	1	0	0	0	1
Proteobacteria	Pseudomonadaceae	369	22	172	9	77	1	2	18	62	1	3
Firmicutes	Amphibacillaceae	203	38	70	0	31	1	3	1	0	8	0
Firmicutes	Bacillaceae	78	6	20	0	6	4	3	4	0	6	5
Firmicutes	Bacillaceae_G	38	3	11	4	20	0	6	4	0	2	4
Proteobacteria	Xanthobacteraceae	230	10	46	4	37	1	3	10	66	9	1

517 **Table 3: Number of unique bacterial homologs (i.e., distinct UHGP-90 protein families) of proteins in**
 518 **a particular xenobiotic class (columns) detected per bacterial family in UHGG (rows).** Highlighted
 519 cells show the top three families per xenobiotic class.

520

521

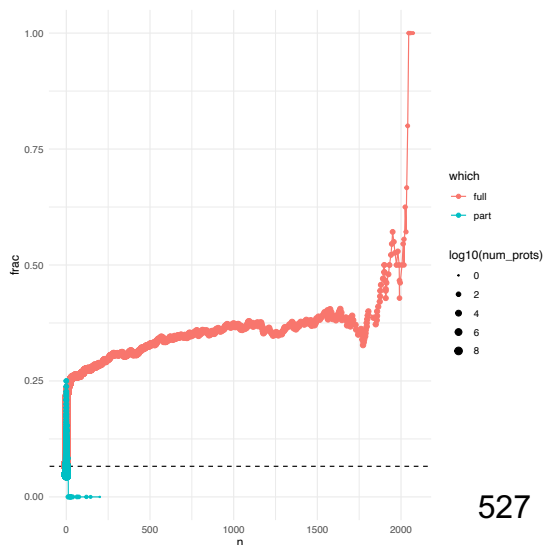
Enzyme	Description	nPart	nFull	From
AOX1	Aldehyde oxidase	46	3	crizotinib, allopurinol, ziprasidone, vortioxetine, aciclovir, pyrazinamide, capmatinib, nicotine iminium ion, thioguanine
XDH	Xanthine dehydrogenase / oxidase	38	0	doxorubicin, 1-methylxanthine, theophylline, allopurinol, pyrazinoic acid, pyrazinamide, mercaptopurine, thioxanthine
DPYD	Dihydropyrimidine dehydrogenase	26	24	fluorouracil
UMPS	Uridine 5'-monophosphate synthase	22	0	fluorouracil

522 **Table 4: Table showing all identified drug-metabolizing enzymes with more partial (nPart) than full-**
523 **length (nFull) gut microbial homologs, together with the drug(s) they metabolize (From).**

524

525

526 Supplementary Tables and Figures



Supplementary Figure 1: Trend in mitochondrial localization for full-length (orange) and split (teal) homologs, as a function of their distribution across gut species. Each dot represents all microbial homologs present in at least a certain number of gut microbial species (x-axis). The size of the dot corresponds to the total number of such microbial homologs. The y-axis shows what fraction of the human homologs are annotated as localizing to the mitochondrion. The overall rate for all proteins is shown by the dashed line.

528

529 **Supplementary Table 1:** GO term enrichment (biological process) for human proteins with more full-length
530 than split homologs. Terms with adjusted p-values below 0.05 are shown.

531 **Supplementary Table 2:** GO term enrichment (biological process) for human proteins with more split than
532 full-length homologs. Terms with adjusted p-values below 0.05 are shown.

533 **Supplementary Table 3:** Drugs metabolized by human proteins with mostly full-length homologs in the gut
534 microbiome.

535 Works Cited

- 536 1. Koppel N, Bisanz JE, Pandelia M-E, Turnbaugh PJ, Balskus EP. Discovery and
537 characterization of a prevalent human gut bacterial enzyme sufficient for the inactivation
538 of a family of plant toxins. Ley RE, editor. *eLife*. 2018;7: e33953. doi:10.7554/eLife.33953
- 539 2. Spanogiannopoulos P, Kyaw TS, Guthrie BGH, Bradley PH, Lee JV, Melamed J, et al.
540 Host and gut bacteria share metabolic pathways for anti-cancer drug metabolism. *Nat*
541 *Microbiol*. 2022;7: 1605–1620. doi:10.1038/s41564-022-01226-5
- 542 3. Dobkin JF, Saha JR, Butler VP, Neu HC, Lindenbaum J. Digoxin-Inactivating Bacteria:
543 Identification in Human Gut Flora. *Science*. 1983;220: 325–327.
- 544 4. Lindenbaum J, Rund DG, Butler VP, Tse-Eng D, Saha JR. Inactivation of digoxin by the
545 gut flora: reversal by antibiotic therapy. *N Engl J Med*. 1981;305: 789–794.
546 doi:10.1056/NEJM198110013051403
- 547 5. Haiser HJ, Gootenberg DB, Chatman K, Sirasani G, Balskus EP, Turnbaugh PJ.
548 Predicting and manipulating cardiac drug inactivation by the human gut bacterium
549 *Eggerthella lenta*. *Science*. 2013;341: 295. doi:10.1126/science.1235872
- 550 6. Malwe AS, Srivastava GN, Sharma VK. GutBug: A Tool for Prediction of Human Gut
551 Bacteria Mediated Biotransformation of Biotic and Xenobiotic Molecules Using Machine
552 Learning. *J Mol Biol*. 2023;435: 168056. doi:10.1016/j.jmb.2023.168056
- 553 7. Guthrie L, Wolfson S, Kelly L. The human gut chemical landscape predicts microbe-
554 mediated biotransformation of foods and drugs. Garrett WS, Turnbaugh P, Turnbaugh P,
555 editors. *eLife*. 2019;8: e42866. doi:10.7554/eLife.42866
- 556 8. Bustion AE, Nayak RR, Agrawal A, Turnbaugh PJ, Pollard KS. SIMMER employs
557 similarity algorithms to accurately identify human gut microbiome species and enzymes
558 capable of known chemical transformations. Redinbo M, Garrett WS, Zimmermann M,
559 editors. *eLife*. 2023;12: e82401. doi:10.7554/eLife.82401
- 560 9. Noda-Garcia L, Tawfik DS. Enzyme evolution in natural products biosynthesis: target- or
561 diversity-oriented? *Curr Opin Chem Biol*. 2020;59: 147–154.
562 doi:10.1016/j.cbpa.2020.05.011
- 563 10. Darby CA, Stolzer M, Ropp PJ, Barker D, Durand D. Xenolog classification.
564 *Bioinformatics*. 2016;33: btw686. doi:10.1093/bioinformatics/btw686
- 565 11. Stolzer M, Siewert K, Lai H, Xu M, Durand D. Event inference in multidomain families
566 with phylogenetic reconciliation. *BMC Bioinformatics*. 2015;16 Suppl 14: S8.
567 doi:10.1186/1471-2105-16-S14-S8
- 568 12. Méheust R, Bhattacharya D, Pathmanathan JS, McInerney JO, Lopez P, Baptiste E.
569 Formation of chimeric genes with essential functions at the origin of eukaryotes. *BMC*
570 *Biol*. 2018;16: 30. doi:10.1186/s12915-018-0500-0

- 571 13. Brueckner J, Martin WF. Bacterial Genes Outnumber Archaeal Genes in Eukaryotic
572 Genomes. *Genome Biol Evol.* 2020;12: 282–292. doi:10.1093/gbe/evaa047
- 573 14. Almeida A, Nayfach S, Boland M, Strozzi F, Beracochea M, Shi ZJ, et al. A unified
574 catalog of 204,938 reference genomes from the human gut microbiome. *Nat Biotechnol.*
575 2021;39: 105–114. doi:10.1038/s41587-020-0603-3
- 576 15. The UniProt Consortium. UniProt: the universal protein knowledgebase in 2021. *Nucleic
577 Acids Res.* 2021;49: D480–D489. doi:10.1093/nar/gkaa1100
- 578 16. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
579 *J Mol Biol.* 1990;215: 403–10. doi:10.1016/S0022-2836(05)80360-2
- 580 17. Aleksander SA, Balhoff J, Carbon S, Cherry JM, Drabkin HJ, Ebert D, et al. The Gene
581 Ontology knowledgebase in 2023. *Genetics.* 2023;224: iyad031.
582 doi:10.1093/genetics/iyad031
- 583 18. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene Ontology:
584 tool for the unification of biology. *Nat Genet.* 2000;25: 25–29. doi:10.1038/75556
- 585 19. Pittis AA, Gabaldón T. Late acquisition of mitochondria by a host with chimaeric
586 prokaryotic ancestry. *Nature.* 2016;531: 101–104. doi:10.1038/nature16941
- 587 20. Abdill RJ, Graham SP, Rubinetti V, Albert FW, Greene CS, Davis S, et al. Integration of
588 168,000 samples reveals global patterns of the human gut microbiome. *bioRxiv.* 2023;
589 2023.10.11.560955. doi:10.1101/2023.10.11.560955
- 590 21. Zong W, Friedman ES, Allu SR, Firman J, Tu V, Daniel SG, et al. Disruption of intestinal
591 oxygen balance in acute colitis alters the gut microbiome. *Gut Microbes.* 2024;16:
592 2361493. doi:10.1080/19490976.2024.2361493
- 593 22. Bradley PH, Pollard KS. Proteobacteria explain significant functional variability in the
594 human gut microbiome. *Microbiome.* 2017;5: 36. doi:10.1186/s40168-017-0244-z
- 595 23. Spinasse LB, Santos AR, Suffys PN, Muxfeldt ES, Salles GF. Different Phenotypes of the
596 NAT2 Gene Influences Hydralazine Antihypertensive Response in Patients with Resistant
597 Hypertension. *Pharmacogenomics.* 2014;15: 169–178. doi:10.2217/pgs.13.202
- 598 24. Butcher NJ, Boukouvala S, Sim E, Minchin RF. Pharmacogenetics of the arylamine N-
599 acetyltransferases. *Pharmacogenomics J.* 2002;2: 30–42. doi:10.1038/sj.tpj.6500053
- 600 25. Payton M, Auty R, Delgoda R, Everett M, Sim E. Cloning and characterization of
601 arylamine N-acetyltransferase genes from *Mycobacterium smegmatis* and
602 *Mycobacterium tuberculosis*: increased expression results in isoniazid resistance. *J
603 Bacteriol.* 1999;181: 1343–1347. doi:10.1128/JB.181.4.1343-1347.1999
- 604 26. Hein DW, Doll MA, Rustan TD, Gray K, Feng Y, Ferguson RJ, et al. Metabolic activation
605 and deactivation of arylamine carcinogens by recombinant human NAT1 and polymorphic

- 606 NAT2 acetyltransferases. *Carcinogenesis*. 1993;14: 1633–1638.
607 doi:10.1093/carcin/14.8.1633
- 608 27. Penning TM, Wangtrakuldee P, Auchus RJ. Structural and Functional Biology of Aldo-
609 Keto Reductase Steroid-Transforming Enzymes. *Endocr Rev*. 2019;40: 447–475.
610 doi:10.1210/er.2018-00089
- 611 28. Kavanagh KL, Jörnvall H, Persson B, Oppermann U. Medium- and short-chain
612 dehydrogenase/reductase gene and protein families: The SDR superfamily: functional
613 and structural diversity within a family of metabolic and regulatory enzymes. *Cell Mol Life*
614 *Sci CMLS*. 2008;65: 3895. doi:10.1007/s00018-008-8588-y
- 615 29. Collins SL, Stine JG, Bisanz JE, Okafor CD, Patterson AD. Bile acids and the gut
616 microbiota: metabolic interactions and impacts on disease. *Nat Rev Microbiol*. 2023;21:
617 236–247. doi:10.1038/s41579-022-00805-x
- 618 30. Qian L, Ouyang H, Gordils-Valentin L, Hong J, Jayaraman A, Zhu X. Identification of Gut
619 Bacterial Enzymes for Keto-Reductive Metabolism of Xenobiotics. *ACS Chem Biol*.
620 2022;17: 1665–1671. doi:10.1021/acscchembio.2c00312
- 621 31. Spanogiannopoulos P, Bess EN, Carmody RN, Turnbaugh PJ. The microbial
622 pharmacists within us: a metagenomic view of xenobiotic metabolism. *Nat Rev Microbiol*.
623 2016;14: 273–287. doi:10.1038/nrmicro.2016.17
- 624 32. Koppel N, Maini Rekdal V, Balskus EP. Chemical transformation of xenobiotics by the
625 human gut microbiota. *Science*. 2017;356: eaag2770. doi:10.1126/science.aag2770
- 626 33. Chen YM, Zhu Y, Lin EC. NAD-linked aldehyde dehydrogenase for aerobic utilization of
627 L-fucose and L-rhamnose by *Escherichia coli*. *J Bacteriol*. 1987;169: 3289.
628 doi:10.1128/jb.169.7.3289-3294.1987
- 629 34. Emadi A, Jones RJ, Brodsky RA. Cyclophosphamide and cancer: golden anniversary.
630 *Nat Rev Clin Oncol*. 2009;6: 638–647. doi:10.1038/nrclinonc.2009.146
- 631 35. Kobayashi Y, Fukami T, Shimizu M, Nakajima M, Yokoi T. Contributions of arylacetamide
632 deacetylase and carboxylesterase 2 to flutamide hydrolysis in human liver. *Drug Metab*
633 *Dispos Biol Fate Chem*. 2012;40: 1080–1084. doi:10.1124/dmd.112.044537
- 634 36. Shimizu M, Fukami T, Ito Y, Kurokawa T, Kariya M, Nakajima M, et al. Indiplon is
635 hydrolyzed by arylacetamide deacetylase in human liver. *Drug Metab Dispos Biol Fate*
636 *Chem*. 2014;42: 751–758. doi:10.1124/dmd.113.056184
- 637 37. Thomsen R, Rasmussen HB, Linnet K, INDICES Consortium. In vitro drug metabolism by
638 human carboxylesterase 1: focus on angiotensin-converting enzyme inhibitors. *Drug*
639 *Metab Dispos Biol Fate Chem*. 2014;42: 126–133. doi:10.1124/dmd.113.053512

- 640 38. Xu G, Zhang W, Ma MK, McLeod HL. Human Carboxylesterase 2 Is Commonly
641 Expressed in Tumor Tissue and Is Correlated with Activation of Irinotecan¹. *Clin Cancer*
642 *Res.* 2002;8: 2605–2611.
- 643 39. Whirl-Carrillo M, McDonagh E, Hebert J, Gong L, Sangkuhl K, Thorn C, et al.
644 Pharmacogenomics Knowledge for Personalized Medicine. *Clin Pharmacol Ther.*
645 2012;92: 414–417. doi:10.1038/clpt.2012.96
- 646 40. Kasahara K, Kerby RL, Zhang Q, Pradhan M, Mehrabian M, Lusic AJ, et al. Gut bacterial
647 metabolism contributes to host global purine homeostasis. *Cell Host Microbe.* 2023;31:
648 1038-1053.e10. doi:10.1016/j.chom.2023.05.011
- 649 41. Yan Y, Wang Z, Zhou Y-L, Gao Z, Ning L, Zhao Y, et al. Commensal bacteria promote
650 azathioprine therapy failure in inflammatory bowel disease via decreasing 6-
651 mercaptopurine bioavailability. *Cell Rep Med.* 2023;4: 101153.
652 doi:10.1016/j.xcrm.2023.101153
- 653 42. Sender R, Fuchs S, Milo R. Revised Estimates for the Number of Human and Bacteria
654 Cells in the Body. *PLOS Biol.* 2016;14: e1002533. doi:10.1371/journal.pbio.1002533
- 655 43. Price G, Patel DA. Drug Bioavailability. StatPearls. Treasure Island (FL): StatPearls
656 Publishing; 2024. Available: <http://www.ncbi.nlm.nih.gov/books/NBK557852/>
- 657 44. Hsu CL, Schnabl B. The gut–liver axis and gut microbiota in health and liver disease. *Nat*
658 *Rev Microbiol.* 2023;21: 719–733. doi:10.1038/s41579-023-00904-3
- 659 45. Steinegger M, Söding J. MMseqs2 enables sensitive protein sequence searching for the
660 analysis of massive data sets. *Nat Biotechnol.* 2017;35: 1026–1028.
661 doi:10.1038/nbt.3988
- 662 46. Corso G, Stärk H, Jing B, Barzilay R, Jaakkola T. DiffDock: Diffusion Steps, Twists, and
663 Turns for Molecular Docking. *arXiv*; 2023. doi:10.48550/arXiv.2210.01776
- 664 47. Abramson J, Adler J, Dunger J, Evans R, Green T, Pritzel A, et al. Accurate structure
665 prediction of biomolecular interactions with AlphaFold 3. *Nature.* 2024;630: 493–500.
666 doi:10.1038/s41586-024-07487-w
- 667 48. Krishna R, Wang J, Ahern W, Sturmfels P, Venkatesh P, Kalvet I, et al. Generalized
668 biomolecular modeling and design with RoseTTAFold All-Atom. *Science.* 2024;384:
669 eadl2528. doi:10.1126/science.adl2528
- 670 49. Maier L, Pruteanu M, Kuhn M, Zeller G, Telzerow A, Anderson EE, et al. Extensive
671 impact of non-antibiotic drugs on human gut bacteria. *Nature.* 2018;555: 623–628.
672 doi:10.1038/nature25979
- 673 50. Timmis JN, Ayliffe MA, Huang CY, Martin W. Endosymbiotic gene transfer: organelle
674 genomes forge eukaryotic chromosomes. *Nat Rev Genet.* 2004;5: 123–135.
675 doi:10.1038/nrg1271

- 676 51. Mondino S, Schmidt S, Buchrieser C. Molecular Mimicry: a Paradigm of Host-Microbe
677 Coevolution Illustrated by Legionella. *mBio*. 2020;11: 10.1128/mbio.01201-20.
678 doi:10.1128/mbio.01201-20
- 679 52. team T pandas development. pandas-dev/pandas: Pandas. Zenodo; 2020.
680 doi:10.5281/zenodo.3509134
- 681 53. Vink R, Gooijer S de, Beedie A, Gorelli ME, Zundert J van, Hulselmans G, et al. pola-
682 rs/polars: Python Polars 0.20.2. Zenodo; 2023. doi:10.5281/zenodo.10413093
- 683 54. R Core Team. R: a language and environment for statistical computing. Vienna, Austria:
684 R Foundation for Statistical Computing; 2024. Available: <https://www.R-project.org/>
- 685 55. Wickham H, Averick M, Bryan J, Chang W, McGowan LD, François R, et al. Welcome to
686 the tidyverse. *J Open Source Softw*. 2019;4: 1686. doi:10.21105/joss.01686
- 687 56. Alexa A, Rahnenfuhrer J. topGO: Enrichment analysis for gene ontology. 2024.
- 688 57. Hochberg Y, Benjamini Y. Controlling the false discovery rate: a practical and powerful
689 approach to multiple testing. *J R Stat Soc Ser B Methodol*. 1995;1: 289–300.
- 690 58. Faith DP. Conservation evaluation and phylogenetic diversity. *Biol Conserv*. 1992;61: 1–
691 10. doi:10.1016/0006-3207(92)91201-3
- 692 59. Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
693 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*.
694 2015;32: 268–274. doi:10.1093/molbev/msu300
- 695 60. Paradis E, Claude J, Strimmer K. APE: Analyses of Phylogenetics and Evolution in R
696 language. *Bioinformatics*. 2004;20: 289–290. doi:10.1093/bioinformatics/btg412
- 697 61. Kembel SW, Cowan PD, Helmus MR, Cornwell WK, Morlon H, Ackerly DD, et al. Picante:
698 R tools for integrating phylogenies and ecology. *Bioinforma Oxf Engl*. 2010;26: 1463–
699 1464.
- 700 62. Yu G, Smith D, Zhu H, Guan Y, Lam TT-Y. ggtree: an R package for visualization and
701 annotation of phylogenetic trees with their covariates and other associated data. *Methods*
702 *Ecol Evol*. 2017;8: 28–36. doi:10.1111/2041-210X.12628
- 703 63. Whirl-Carrillo M, Huddart R, Gong L, Sangkuhl K, Thorn CF, Whaley R, et al. An
704 Evidence-Based Framework for Evaluating Pharmacogenomics Knowledge for
705 Personalized Medicine. *Clin Pharmacol Ther*. 2021;110: 563–572. doi:10.1002/cpt.2350
- 706 64. Seal RL, Braschi B, Gray K, Jones TEM, Tweedie S, Haim-Vilmovsky L, et al.
707 Genenames.org: the HGNC resources in 2023. *Nucleic Acids Res*. 2023;51: D1003–
708 D1009. doi:10.1093/nar/gkac888

- 709 65. Sayers EW, Cavanaugh M, Clark K, Pruitt KD, Schoch CL, Sherry ST, et al. GenBank.
710 Nucleic Acids Res. 2020;49: D92–D96. doi:10.1093/nar/gkaa1023
- 711 66. Cantelli G, Bateman A, Brooksbank C, Petrov AI, Malik-Sheriff RS, Ide-Smith M, et al.
712 The European Bioinformatics Institute (EMBL-EBI) in 2021. Nucleic Acids Res. 2022;50:
713 D11–D19. doi:10.1093/nar/gkab1127
- 714