

## Supplementary information

### Amino acid auxotrophies in human gut bacteria are linked to higher microbiome diversity and long-term stability

Svenja Starke<sup>1</sup>, Danielle MM Harris<sup>1,2</sup>, Johannes Zimmermann<sup>3,4</sup>, Sven Schuchardt<sup>5</sup>, Mhmd Oumari<sup>2</sup>, Derk Frank<sup>6,7</sup>, Corinna Bang<sup>2</sup>, Philip Rosenstiel<sup>2</sup>, Stefan Schreiber<sup>2,8</sup>, Norbert Frey<sup>6,7,#</sup>, Andre Franke<sup>2</sup>, Konrad Aden<sup>2,8,\*</sup>, Silvio Waschina<sup>1,\*</sup>

<sup>1</sup> Institute of Human Nutrition and Food Science, Nutriinformatics, Kiel University, Kiel, Germany

<sup>2</sup> Institute of Clinical Molecular Biology, Kiel University, Kiel, Germany

<sup>3</sup> Research Group Evolutionary Ecology and Genetics, Zoological Institute, University of Kiel, Kiel, Germany

<sup>4</sup> Max Planck Institute for Evolutionary Biology, Plön, Germany

<sup>5</sup> Fraunhofer Institute for Toxicology and Experimental Medicine (ITEM), Hanover, Germany

<sup>6</sup> Department of Internal Medicine III, University Medical Center Schleswig-Holstein, Kiel, Germany

<sup>7</sup> German Centre for Cardiovascular Research (DZHK), Partner site Hamburg, Kiel, Lübeck, Germany

<sup>8</sup> Department of Internal Medicine I, University Medical Center Schleswig-Holstein, Kiel, Germany

# Current Affiliation: Department of Internal Medicine III, University Hospital Heidelberg, Heidelberg, Germany

\*Correspondence: [s.waschina@nutrinf.uni-kiel.de](mailto:s.waschina@nutrinf.uni-kiel.de) ; [k.aden@ikmb.uni-kiel.de](mailto:k.aden@ikmb.uni-kiel.de)

## 27 Supplementary Material and Methods

28

### 29 **Reconstruction of genome-scale metabolic models**

30 Genome-scale metabolic models of prokaryotic genomes were reconstructed using

31 *gapseq*(1). In brief, the *gapseq* reconstruction workflow consisted of five steps: (i) Reaction

32 and pathway prediction, (ii) prediction of metabolite cross-membrane transporters, (iii)

33 reconstruction of a draft metabolic network based on the results from *i* and *ii*, (iv)

34 estimation of an organism-specific growth medium-based on the predicted metabolic

35 capabilities, and (v) gap filling of the metabolic network to enable biomass production using

36 flux balance analysis. Model reconstructions were limited to bacterial genomes marked as

37 representative species in the HRGM collection. Further, genomes with an estimated

38 contamination percentage of  $\leq 2\%$  or a completion  $\geq 85\%$  were included. Based on these

39 filters, 3 687 bacterial genomes were subject to metabolic model reconstruction. Among

40 those genomes, 22% are from bacterial isolates, and 78% are metagenome-assembled

41 genomes.

42 A recent computational study has shown in a systematic analysis of isolate- and

43 metagenome-assembled genomes that the gap-filling medium strongly impacts the

44 auxotrophies predicted by genome-scale metabolic modelling(2). We note that with the

45 reconstruction procedure used in this study, we do not rely on an arbitrarily defined gap-

46 filling medium, which is used for every metabolic network model. Instead, for each draft

47 network, a genome-specific gap-filling medium is predicted (see section "Prediction of a

48 genome-specific gap-filling medium" below for details). In brief, if the medium prediction

49 algorithm of *gapseq* finds a known biosynthetic pathway for a specific amino acid in the

50 draft metabolic network, the amino acid will not be part of the resulting predicted medium,

51 as the compound is likely not required from the growth environment, thus, also not  
52 necessary for subsequent gap filling. In contrast, if the medium prediction algorithm does  
53 not detect at least one complete known biosynthetic pathway for a specific amino acid, the  
54 respective compound is added to the outcome gap-filling medium based on the rationale  
55 that this amino is a putative essential compound that needs to be obtained from the growth  
56 environment. However, it is important to note that this case does not directly imply that the  
57 organism is auxotrophic for the specific amino acid, as the subsequent gap-filling algorithm  
58 might add reactions to the model that complete a biosynthesis route from other  
59 compounds in the gap-filling medium to the amino acid. Such reactions are only added in  
60 cases where a gene was found in the query genome that displays sequence similarity to a  
61 reference gene sequence with the respective enzymatic function but where the sequence  
62 similarity was not high enough to pass the threshold of the bitscore 200 to be directly  
63 included in the draft network. For details on the gap-filling algorithm implemented in  
64 *gapseq*, please refer to the original *gapseq* publication(1).

65 Taken together, the model reconstruction and auxotrophy prediction that we used for the  
66 present study do not depend on one gap-filling medium composition that is defined for all  
67 organisms but adjusts the medium for each organism based on its genome information and  
68 by using the multi-step gap-filling algorithm that is implemented in *gapseq*.

69

#### 70 **Prediction of a genome-specific gap-filling medium**

71 The genome-scale metabolic network reconstruction process of *gapseq* requires a gap-filling  
72 medium for the final gap-filling step (see above). Here, we used a medium prediction  
73 feature (module "*gapseq medium*") of the *gapseq* software, which can be plugged into the  
74 reconstruction workflow between the homology-based generation of a draft metabolic

75 network and the gap-filling algorithm. We provided the additional command line option “-c  
76 cpd00007:0” to ensure that the predicted medium does not contain oxygen (compound  
77 identifier: cpd00007). The algorithm for medium prediction tests which pathways and  
78 reactions are absent or present in the draft metabolic model. Whether a particular  
79 compound is added to the medium is decided using logical expressions that include  
80 variables for the presence (TRUE) and absence (FALSE) of pathways and reactions within the  
81 design network (see Supplementary Table S6 for all compounds and their logical  
82 expressions). For example, the disaccharide lactose (ModelSEED ID: cpd00208), has the  
83 logical expression ("LACTOSECAT-PWY" | "LACTOSEUTIL-PWY" | "BGALACT-PWY"), which  
84 means, that lactose is added to the medium if one of three known lactose degradation  
85 pathways as defined in MetaCyc(3) is already present in the draft network. In particular, for  
86 amino acids, the medium prediction module uses a similar approach to the auxotrophy  
87 prediction tool GapMind(4), which tests if a known biosynthetic pathway for a specific  
88 amino acids exists based on sequence homology. The amino acid biosynthetic pathways that  
89 are considered are also those that are defined in MetaCyc(3). For instance, if none of the  
90 five known biosynthesis pathways in prokaryotes for lysine  
91 (<https://metacyc.org/META/NEW-IMAGE?type=PATHWAY&object=LYSINE-SYN>) is found,  
92 lysine is added to the gap-filling medium.

93 The medium prediction module of gapseq considered 74 compounds (Supplementary Table  
94 S6), including inorganic compounds, carbohydrates, amino acids, other carboxylic acids, and  
95 vitamins. Most of those potential nutrients are also compounds that can be found in the  
96 growth environment of colonic microorganisms, such as fibers (e.g., pectin, inulin), other  
97 dietary compounds (e.g., sulfoquinovose, daidzein), constituents of the mucins (e.g., *N*-  
98 acetylneuraminate, *N*-acetylneuraminate) and inorganics (e.g., H<sub>2</sub>, H<sub>2</sub>S, H<sub>2</sub>O).

99 Besides the enumeration of available nutrients, a gap-filling medium requires their  
100 individual maximum uptake rates by the microorganism. The medium prediction  
101 implemented in *gapseq* uses rates commonly used in manually curated genome-scale  
102 metabolic network models, e.g.,  $0.1 \text{ mmol} \cdot \text{gDW}^{-1} \cdot \text{hr}^{-1}$  for amino acids or  $5 \text{ mmol} \cdot \text{gDW}^{-1} \cdot \text{hr}^{-1}$   
103 for monosaccharides. In the case of oligo- and polysaccharides, the maximum uptake rates  
104 are scaled to allow the same uptake rate per subunit (e.g.,  $5 \text{ mmol} \cdot \text{gDW}^{-1} \cdot \text{hr}^{-1}$  for the  
105 monosaccharide D-glucose and  $2.5 \text{ mmol} \cdot \text{gDW}^{-1} \cdot \text{hr}^{-1}$  for the disaccharide maltose).

106

### 107 **Validation of auxotrophy predictions**

108 To validate our *in silico* predicted auxotrophies, we collected the genome sequences from  
109 NCBI RefSeq for 36 bacterial strains for which experimental data were available on amino  
110 acid auxotrophies/prototrophies. The majority of these strains were already summarized  
111 previously(5). In that study, the authors even list more than the 36 strains we analyzed here.  
112 This is because we excluded cases where we could not find a genome assembly of the exact  
113 strain used in the referenced experimental study. Moreover, we excluded the entries for  
114 species belonging to the genus *Bifidobacterium* since their auxotrophy for cysteine is  
115 ambiguous: Some studies report cysteine as an essential nutrient for growth(6,7), while  
116 genomic analysis and genome-scale metabolic modeling indicated the presence of cysteine  
117 biosynthetic pathways(8,9). This ambiguity most likely stems from the fact that cysteine  
118 biosynthesis in *Bifidobacterium* species depends on the available sulfur source(8). Genomic  
119 analysis of *Bifidobacterium bifidum* PRL2010, for instance, suggested the strain's inability to  
120 use sulfate as a sulfur source, while hydrogen sulfide or methionine could potentially serve  
121 as a sulfur source for the biosynthesis of cysteine(6,8).

122 Genome-scale metabolic models for all 36 strains were reconstructed as described above.  
123 Auxotrophies were predicted with the method described in the main manuscript.  
124 In addition to the *gapseq*-reconstructed models, we also predicted auxotrophies for 20 of  
125 the 36 bacterial strains using genome-scale models from the AGORA2 collection(10)  
126 (Supplementary Table S2). Auxotrophies were predicted in the same manner as for *gapseq*  
127 models. In contrast to the *gapseq* models, which can contain only free amino acids and not  
128 peptides in the predicted medium, some AGORA2 models have exchange reactions with  
129 lower bounds  $< 0$  for dipeptides. In those cases, and for predicting the auxotrophy status for  
130 amino acid  $x$ , we changed the lower bound to 0 for all exchange reactions of dipeptides  
131 involving amino acid  $x$ . At the same time, we introduced a new inflow reaction of the non- $x$   
132 amino acid moiety to the model to predict only the essentiality of  $x$  and not of the other  
133 amino acids in the respective peptides. The sensitivity, specificity, and accuracy of  
134 auxotrophy predictions were calculated for *gapseq* models ( $n=36$ ) and AGORA2 models  
135 ( $n=20$ ).  
136 As an additional auxotrophy prediction validation step, 124 genome-scale metabolic models  
137 were reconstructed for bacterial strains that were reported by Price, 2023, to be able to  
138 grow in a defined growth medium containing no amino acids(11). Thus, these 124 organisms  
139 are known amino acid prototrophs and can be used to estimate the rate of false auxotrophy  
140 predictions. The original publication by Price reported 127 genomes of prototrophs;  
141 however, 3 of the corresponding genome assemblies (GCF\_000014265.1,  
142 GCF\_000020545.1, GCF\_900188395.1) were suppressed on RefSeq at the time we  
143 performed the analysis in March 2023. Auxotrophies for the 124 genome assemblies of this  
144 prototroph collection were predicted as described above, and results are summarized in  
145 Supplementary Table S3.

146

## 147 **Metagenome data processing**

148 Metagenomic reads were subject to quality control and filtering using the 'qc' workflow  
149 from the metagenome-atlas pipeline tool v2.9.0(12). In detail, reads were (i) deduplicated,  
150 (ii) quality filtered, and (iii) decontaminated. Modules from the BBmap suite v37.99 (BBMap  
151 - Bushnell B. - [sourceforge.net/projects/bbmap/](https://sourceforge.net/projects/bbmap/)) were used for all three steps. In the  
152 deduplication step (i), the BBmap module *clumpify.sh* was used with the parameters  
153 "dedupe=t dupesubs=2", which removed duplicate reads with a maximum of 2 substitutions  
154 between duplicates. The quality filter (ii) employed the BBmap module *bbduk.sh* with the  
155 parameters "hdist=1 ktrim=r mink=8 trimq=10 qtrim=rl minlength=51 maxns=-1  
156 minbasefrequency=0.05" and otherwise default options. This quality filter trimmed reads  
157 from the right if adapter sequences were detected, trimmed reads on both sides from the  
158 first base with a quality score below 10, removed sequences that were shorter than 51 bp  
159 after trimming, removed sequences with ambiguous base calls (i.e., "N"s), and removed  
160 reads if any base had a frequency of less than 5%. Finally, reads that are likely  
161 contaminations from the human host genome or Illumina PhiX sequences were removed  
162 using the BBmap module *bbsplit.sh* using the option "maxratio=0.65" and otherwise default  
163 parametrization. This tool tested if specific reads mapped to the host genome or PhiX  
164 sequences based on sequence similarity and mapped reads were discarded from the  
165 sample's fastq files. For the decontamination step (iii), the human reference genome  
166 assembly 'Genome Reference Consortium Human Build 38' (GRCh38) was used, in which  
167 low entropy regions (entropy < 0.7) were masked using the *bbmask.sh* tool within the  
168 BBmap suite. Moreover, regions that display high similarity to prokaryotic rRNA genes were  
169 additionally masked. To this end, prokaryotic small and large subunit rRNA gene sequences

170 were retrieved from SILVA version 138.1(13) and shredded into shorter (80 bp) sequences  
171 with 40 bp overlaps using *shred.sh*. Shredded sequences were aligned to GRCh38 with a  
172 minimum identity of 85% and maximum indel length of 2 bp. Regions in GRCh38 with  
173 alignment hits were masked.

174 As mentioned in the main manuscript, we used the Human Reference Gut Microbiome  
175 'HRGM' catalog(14) as reference genomes for quantifying representative microbial genomes  
176 in the metagenomic data sets. We had chosen this collection, as it was the latest published  
177 human gut microorganism genome collection when we were finalizing the results of the  
178 present study. Furthermore, the HRGM collection contains 780 species-level representative  
179 genomes, which were absent in previous genome collections and assembled from  
180 metagenome samples from before under-represented Asian countries, namely Korea,  
181 Japan, and India. For each metagenome sample, the relative abundance of HRGM genomes  
182 was estimated using coverM(15) v0.6.1 with default parametrization of the module `coverm  
183 genome`.

184

#### 185 **Targeted metabolomics of blood samples**

186 Serum samples were collected using serum s-monovette (9ml, Sarstedt, Germany). Samples  
187 were incubated upright at RT for 30 min. and centrifuged (10 min., 2000 x g). Serum was  
188 aliquoted in 500µl tubes and stored at -80°C. Metabolite quantification for serum was  
189 performed by liquid chromatography tandem mass spectrometry (LC-MS-MS) using the MxP  
190 Quant 500 kit (Biocrates Life Sciences AG, Innsbruck, Austria) according to the  
191 manufacturer's instructions. The MxP Quant 500 kit simultaneously measures 630  
192 metabolites covering 14 small molecule and 12 different lipid classes. It combines flow  
193 injection analysis tandem mass spectrometry (FIA-MS/MS) using SCIEX 5500 QTrap mass



194 spectrometer (SCIEX, Darmstadt, Germany) for lipids and liquid chromatography tandem  
195 mass spectrometry (LC-MS/MS) using Agilent 1290 Infinity II liquid chromatography (Santa  
196 Clara, CA, USA) coupled with a SCIEX 5500 QTrap mass spectrometer for small molecules  
197 using multiple reaction monitoring (MRM) to detect the analytes. Data evaluation for serum  
198 metabolite concentrations and quality assessment was performed with the software SCIEX  
199 Analyst software (Version 1.7.2) and the MetIDQ™ software package (Oxygen-DB110-3023),  
200 which is an integral part of the MxP Quant 500 kit.

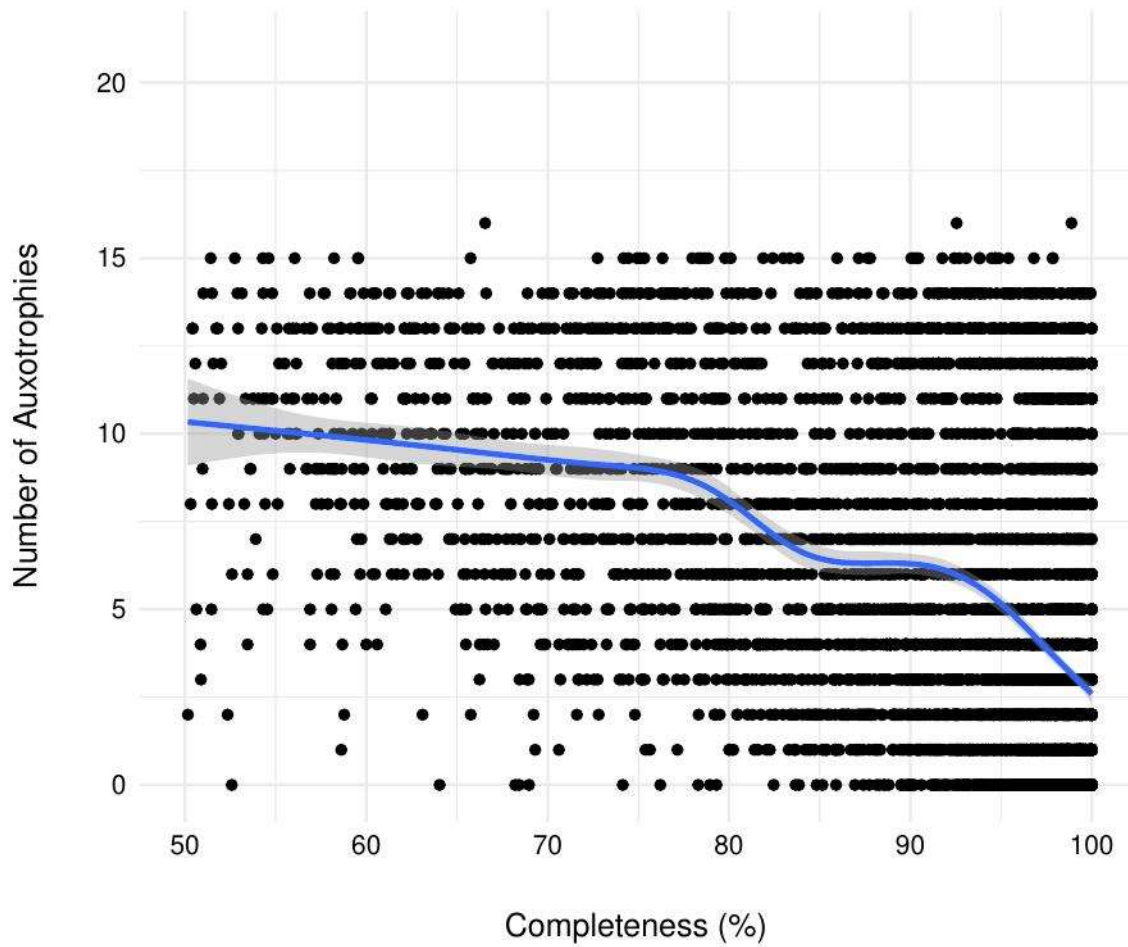
201 For downstream statistical analysis, the serum metabolome data were pre-processed by  
202 imputing missing specific values using a random forests approach as implemented in the R-  
203 package 'missForest' and the function with the same name in default parametrization (16).  
204 This imputation was limited to missing values for metabolites, which have less than 20%  
205 missing values across the data set.

206 With a partial Spearman correlation, the association between the frequency of auxotrophic  
207 bacteria and serum metabolites and other hematology parameters from the DZHK cohort  
208 was evaluated (17). We adjusted for the potential confounders sex, age, and BMI. *p* values  
209 were corrected for multiple testing using the False Discovery Rate (FDR) method.

210

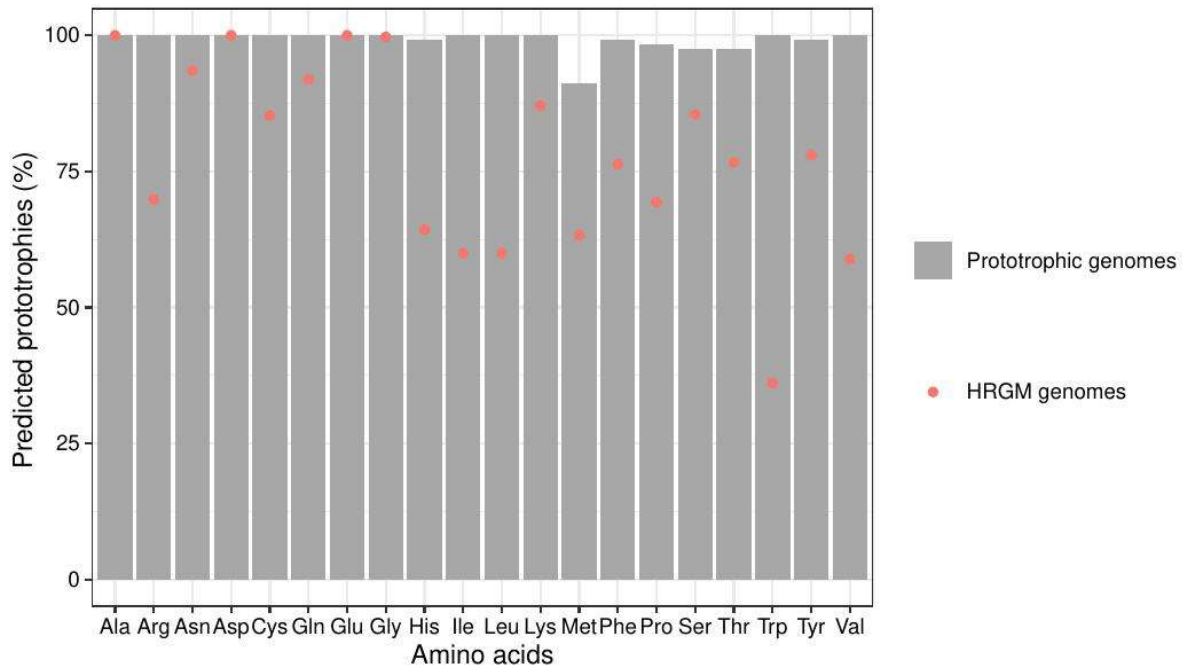
211 Supplementary Figures

212



213

214 **Supp. Figure S1:** Estimated genome completeness and the predicted number of  
215 auxotrophies for 3 687 genomes of representative species from the Human Reference Gut  
216 Microbiome (HRGM) collection(14). The blue line shows the regression line ( $r^2 = -0.50$ ,  $p \leq$   
217  $2.2e-16$ ).



218

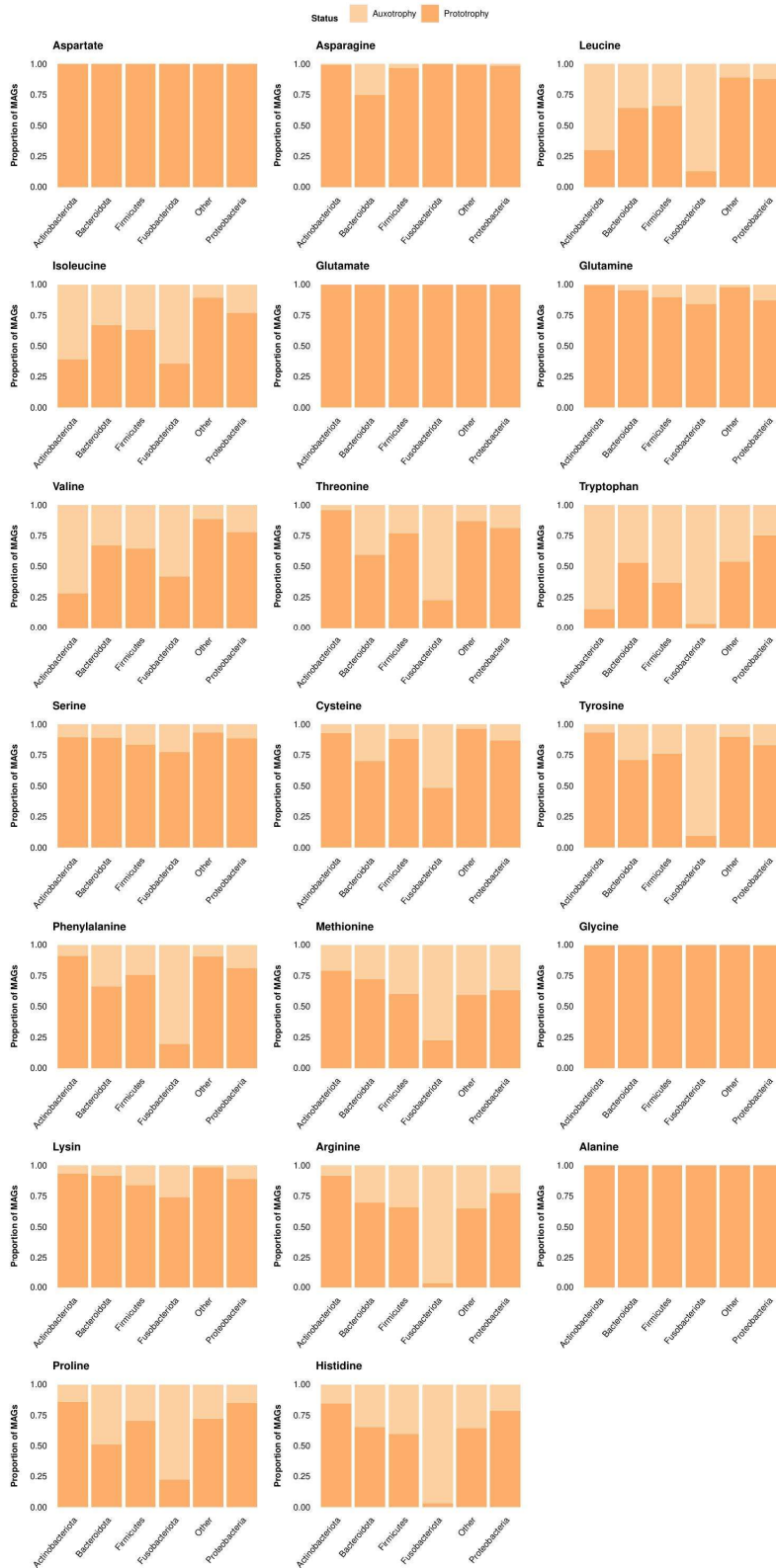
219 **Supp. Figure S2:** Percentage of *in silico* predicted prototrophies with metabolic modeling in

220 124 genomes known to be prototrophic(11) from laboratory experiments (grey bars). The

221 red dots indicate the frequency of prototrophies among 3 687 genomes from human gut

222 bacteria(14).

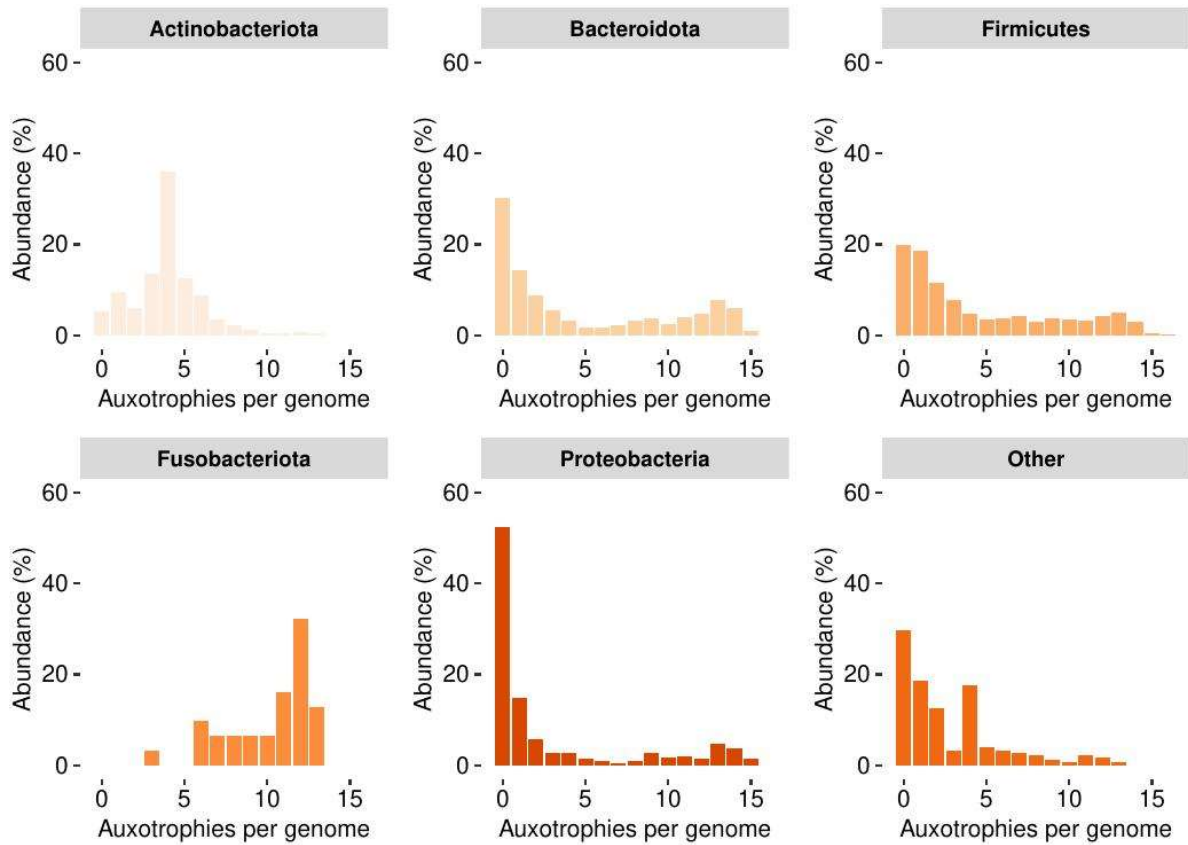
223



224

225 **Supp. Figure S3:** Overview of the proportions of auxotrophy to prototrophy genomes per

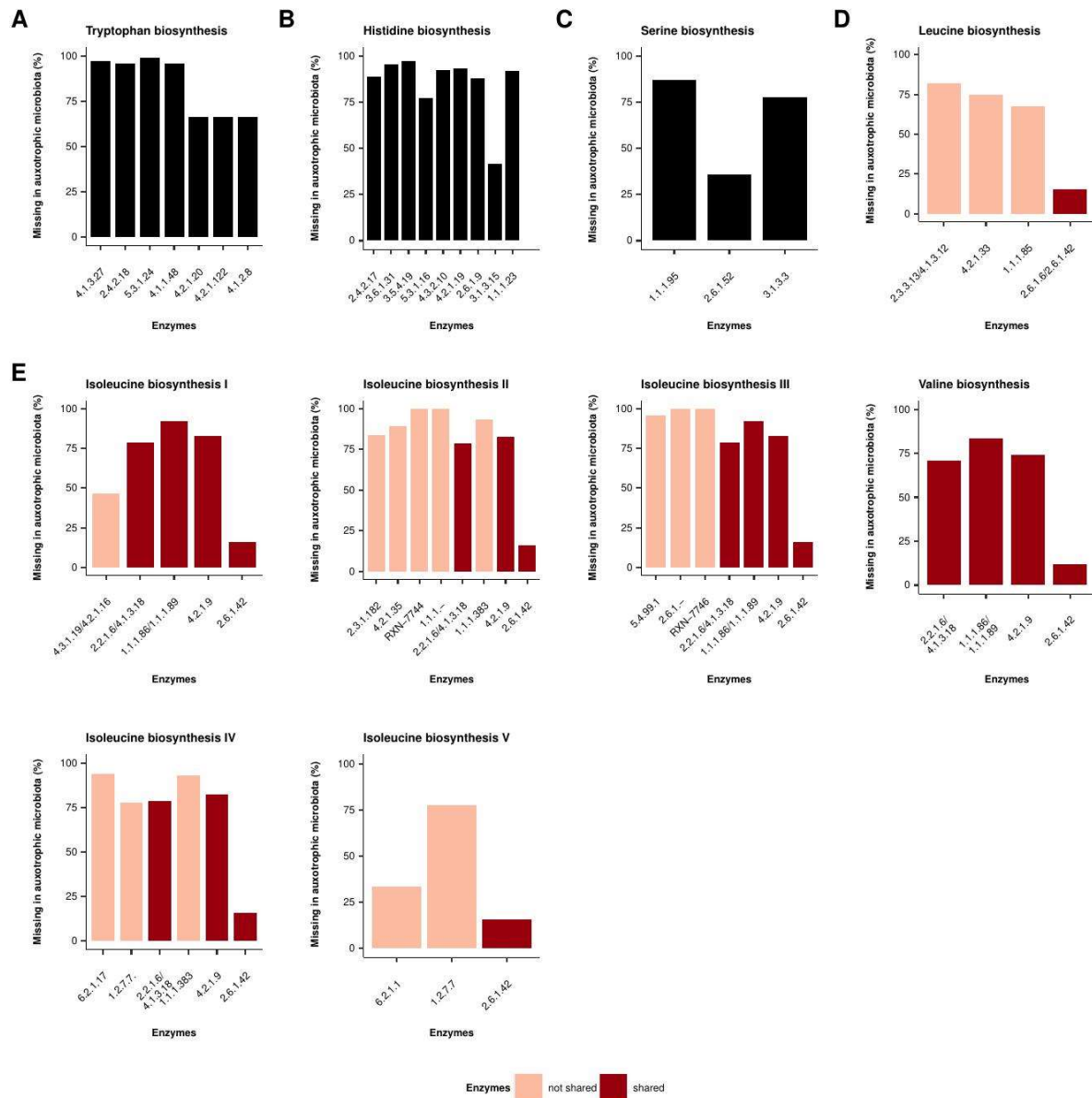
226 phylum from the HRGM catalog(14).



227

228 **Supp. Figure S4:** Number of auxotrophies for every phylum. Other is a category that  
 229 combines different phyla with a lower abundance in the overall HRGM catalogue(14) and for  
 230 a reduction of complexity.

231



232

233 **Supp. Figure S5:** The bar plots display the abundance of missing enzymes in the biosynthesis

234 pathways for several amino acids and BCAA biosynthesis pathways in auxotrophic bacteria.

235 The order of the enzymes in the bar plots represents the one in the pathway. (A) Missing

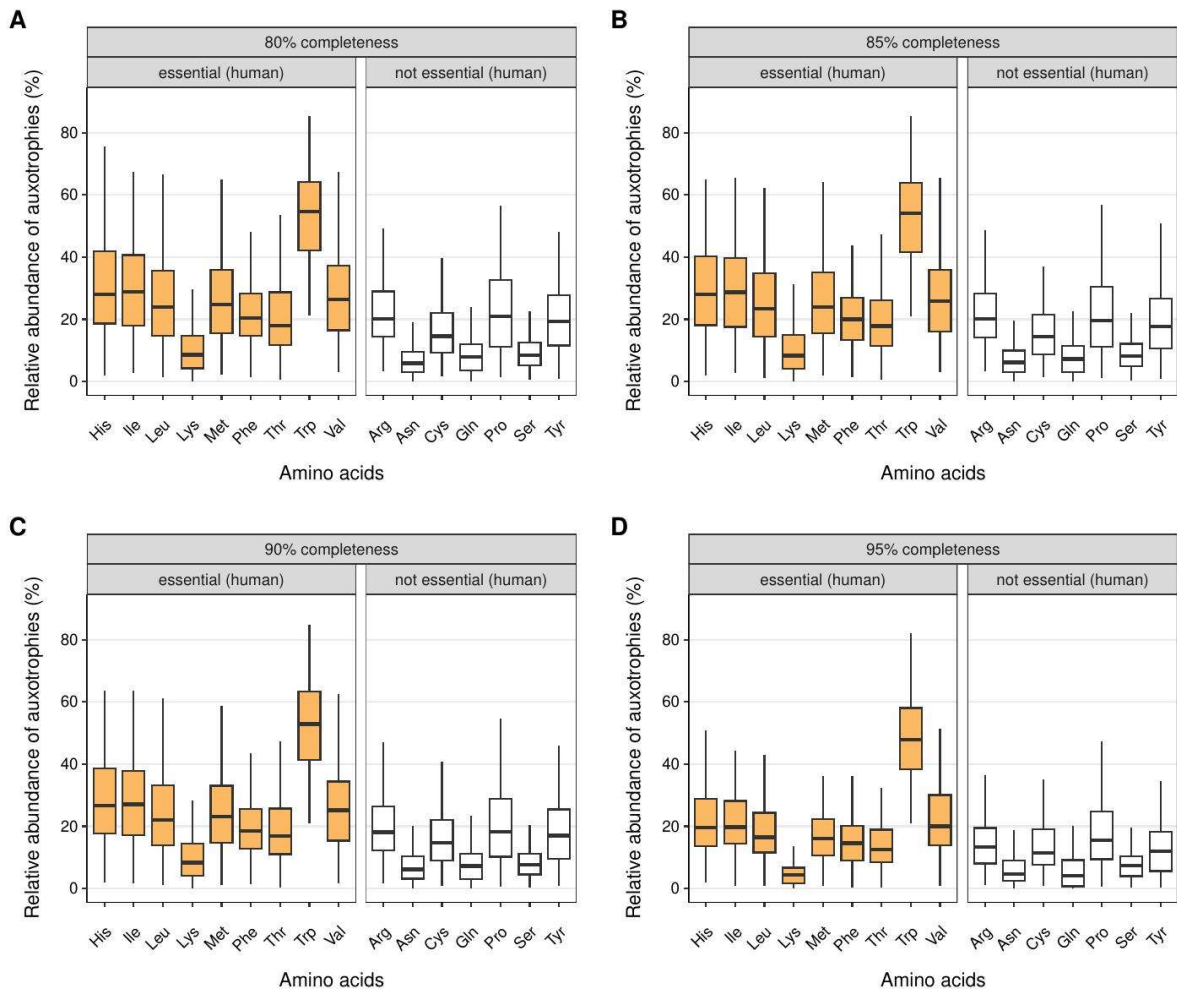
236 enzymes in the tryptophan pathway of tryptophan auxotrophic bacteria, (B) Missing

237 enzymes in the histidine pathway of histidine auxotrophic bacteria, (C) Missing enzymes in

238 the chorismate pathway of chorismate auxotrophic bacteria, (D) Missing enzymes in the

239 serine pathway of serine auxotrophic bacteria, (E) Comparison of the BCAA pathways of

240 isoleucine, leucine, and valine auxotrophic bacteria, the colors indicate which enzymes are  
241 shared in the biosynthesis pathways, the definition of the pathways are based on MetaCyc.  
242



243

244 **Supp. Figure S6:** Relative abundance of predicted amino acids auxotrophs depending on the

245 completeness cutoff for reference genomes from the HRGM catalog. Four different genome

246 completeness cutoffs were tested: 80% (A), 85% (B, same as Figure 4A), 90% (C), and 95%

247 (D). The results indicate that the distribution of the relative abundance of predicted amino

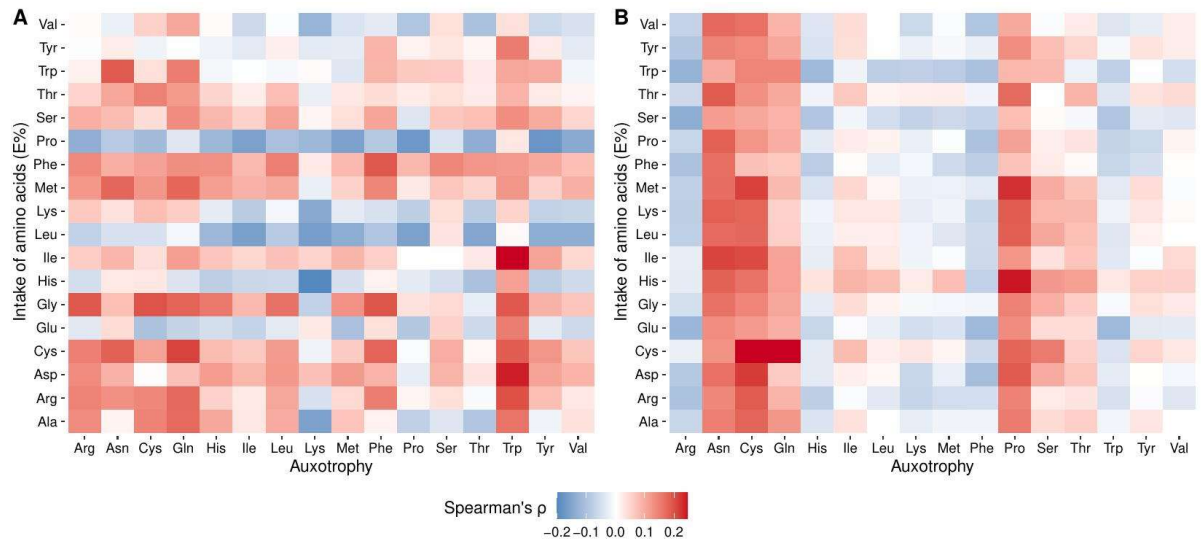
248 acid auxotrophies was stable with respect to the chosen completeness cut-off for reference

249 genome filtering.

250

251





252

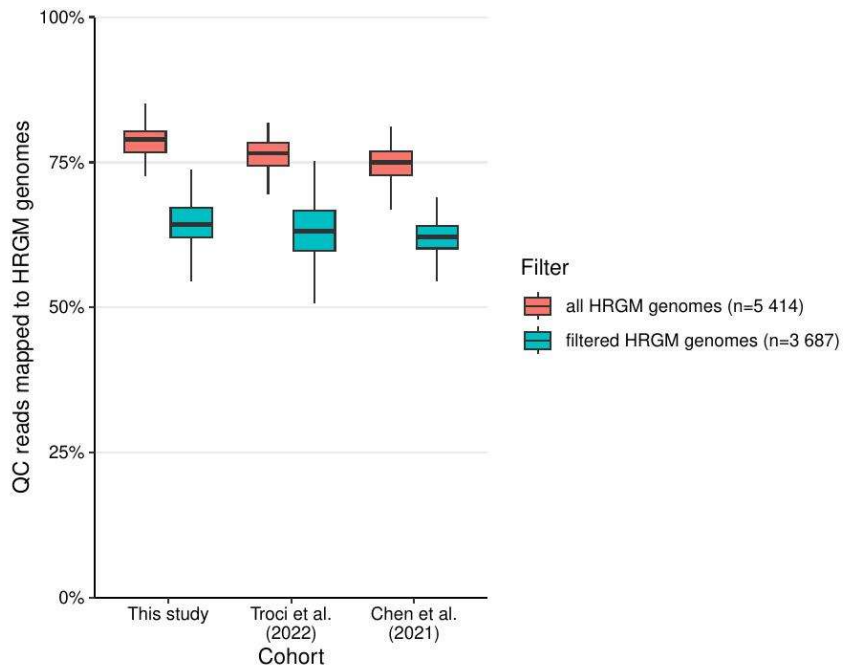
253 **Supp. Figure S7:** Spearman correlation between the dietary intake of amino acids and the

254 frequency of amino acid auxotrophic bacteria in the gut microbiomes, (A) at the beginning

255 of the study, (B) at the end of the study (3 years later). No statistically significant

256 associations were found (FDR-corrected  $p$  value  $>0.05$ ).

257



258

259 **Supp. Figure S8:** Percentage of quality-controlled metagenomic reads from three cohorts

260 (this study, Troci *et al.* 2022 (18), and Chen *et al.* 2021 (19)) to reference genomes from the

261 HRGM catalog.

## 262 References

- 263 1. Zimmermann J, Kaleta C, Waschina S. gapseq: informed prediction of bacterial  
264 metabolic pathways and reconstruction of accurate metabolic models. *Genome Biol.*  
265 2021 Dec;22(1):81.
- 266 2. Borer B, Magnúsdóttir S. The media composition as a crucial element in high-throughput  
267 metabolic network reconstruction. *Interface Focus.* 2023 Apr 6;13(2):20220070.
- 268 3. Caspi R, Billington R, Fulcher CA, Keseler IM, Kothari A, Krummenacker M, et al. The  
269 MetaCyc database of metabolic pathways and enzymes. *Nucleic Acids Res.* 2018 Jan  
270 4;46(D1):D633–9.
- 271 4. Price MN, Deutschbauer AM, Arkin AP. GapMind: Automated Annotation of Amino Acid  
272 Biosynthesis. Hallam SJ, editor. *mSystems.* 2020 Jun 30;5(3):e00291-20.
- 273 5. Ashniev GA, Petrov SN, Iablokov SN, Rodionov DA. Genomics-Based Reconstruction and  
274 Predictive Profiling of Amino Acid Biosynthesis in the Human Gut Microbiome.  
275 *Microorganisms.* 2022 Apr;10(4):740.
- 276 6. Ferrario C, Duranti S, Milani C, Mancabelli L, Lugli GA, Turrone F, et al. Exploring Amino  
277 Acid Auxotrophy in *Bifidobacterium bifidum* PRL2010. *Front Microbiology.* 2015  
278 Nov;6:1331.
- 279 7. Veda M, Nakamoto S, Nakai R, Takagi A. Establishment of a defined minimal medium  
280 and isolation of auxotrophic mutants for *Bifidobacterium bifidum* Es 5. *J Gen Appl*  
281 *Microbiol.* 1983;29(2):103–14.
- 282 8. Schöpping M, Gaspar P, Neves AR, Franzén CJ, Zeidan AA. Identifying the essential  
283 nutritional requirements of the probiotic bacteria *Bifidobacterium animalis* and  
284 *Bifidobacterium longum* through genome-scale modeling. *npj Syst Biol Appl.* 2021 Dec  
285 9;7(1):47.
- 286 9. Lee JH, O’Sullivan DJ. Genomic Insights into Bifidobacteria. *Microbiol Mol Biol Rev.* 2010  
287 Sep;74(3):378–416.
- 288 10. Heinken A, Hertel J, Acharya G, Ravcheev DA, Nyga M, Okpala OE, et al. Genome-scale  
289 metabolic reconstruction of 7,302 human microorganisms for personalized medicine.  
290 *Nat Biotechnol.* 2023 Sept;41:1320-1331.
- 291 11. Price M. Erroneous predictions of auxotrophies by CarveMe. *Nat Ecol Evol.* 2023  
292 Feb;7(2):194–5.
- 293 12. Kieser S, Brown J, Zdobnov EM, Trajkovski M, McCue LA. ATLAS: a Snakemake workflow  
294 for assembly, annotation, and genomic binning of metagenome sequence data. *BMC*  
295 *Bioinformatics.* 2020 Dec;21(1):257.
- 296 13. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, et al. The SILVA ribosomal  
297 RNA gene database project: improved data processing and web-based tools. *Nucleic*  
298 *Acids Res.* 2012 Nov 27;41(D1):D590–6.

- 299 14. Kim CY, Lee M, Yang S, Kim K, Yong D, Kim HR, et al. Human reference gut microbiome  
300 catalog including newly assembled genomes from under-represented Asian  
301 metagenomes. *Genome Med.* 2021 Dec;13(1):134.
- 302 15. Woodcroft BJ. CoverM [Internet]. 2023 [cited 2023 Aug 20]. Available from:  
303 <https://github.com/wwood/CoverM>
- 304 16. Daniel J. Stekhoven. missForest: Nonparametric Missing Value Imputation using Random  
305 Forest. *Bioinformatics.* 2012;28(1):112–8.
- 306 17. Liu Q, Li C, Wanga V, Shepherd BE. Covariate-adjusted Spearman’s rank correlation with  
307 probability-scale residuals: Covariate-Adjusted Spearman’s Rank Correlation. *Biometrics.*  
308 2018 Jun;74(2):595–605.
- 309 18. Troci A, Rausch P, Waschina S, Lieb W, Franke A, Bang C. Long-Term Dietary Effects on  
310 Human Gut Microbiota Composition Employing Shotgun Metagenomics Data Analysis.  
311 *Mol Nutr Food Res.* 2022 Jun 27;e2101098.
- 312 19. Chen L, Wang D, Garmaeva S, Kurilshikov A, Vich Vila A, Gacesa R, et al. The long-term  
313 genetic stability and individual specificity of the human gut microbiome. *Cell.* 2021  
314 Apr;184(9):2302–2315.e12.
- 315