# Human reference gut microbiome comprising 5,414 prokaryotic species, including newly assembled genomes from under-represented Asian metagenomes

Chan Yeong Kim[1†], Muyoung Lee[1†], Sunmo Yang[1], Kyungnam Kim[2], Dongeun Yong[2], Hye Ryun Kim[3], and Insuk Lee[1*]

[1] Department of Biotechnology, College of Life Science & Biotechnology, Yonsei University, Seoul 03722, Korea

[2] Department of Laboratory Medicine, Research Institute of Bacterial Resistance, College of Medicine, Yonsei University, Seoul 03722, Korea

[3] Division of Medical Oncology, Department of Internal Medicine, Yonsei Cancer Center, College of Medicine, Yonsei University, Seoul 03722, Korea

[†] These authors contributed equally to this work

[*] Corresponding author:

Insuk Lee

Tel: +82-10-4186-8706, E-mail: insuklee@yonsei.ac.kr

Short title: Human reference gut microbiome

Key words: metagenomic shotgun sequencing, human gut microbiome, metagenome-assembled genome

## Abstract

Metagenome sampling bias for geographical location and lifestyle is partially responsible for the incomplete catalog of reference genomes of gut microbial species. Here, we present a substantially expanded microbiome catalog, the Human Reference Gut Microbiome (HRGM). Incorporating newly assembled 29,082 genomes from 845 fecal samples collected from three under-represented Asian countries—Korea, India, and Japan—the HRGM contains 232,098 non-redundant genomes of 5,414 representative prokaryotic species, >103 million unique proteins, and >274 million single-nucleotide variants. This is an over 10% increase from the largest reference database. The newly assembled genomes were enriched for members of the *Bacteroidaceae* family, including species associated with high-fiber and seaweed-rich diet. Single-nucleotide variant density was positively associated with the speciation rate of gut commensals. Ultra-deep sequencing facilitated the assembly of genomes of low-abundance taxa, and deep sequencing (>20 million read pairs) was needed for the profiling of low-abundance taxa. Importantly, the HRGM greatly improved the taxonomic and functional classification of sequencing reads from fecal samples. Finally, mapping homologous sequences for human auto-antigens onto the HRGM genomes revealed the association of commensal bacteria with high cross-reactivity potential with autoimmunity. The HRGM (www.mbiomenet.org/HRGM/) will facilitate the identification and functional analysis of disease-associated gut microbiota.

## Introduction

Human gut microbiome is considered the "second human genome" and plays a crucial role in various diseases[1,2]. Therefore, targeting gut microbes and their functional elements may provide novel therapeutic opportunities. The assembly of human reference genome, together with a catalog of protein-coding genes and genomic variants, led us to the era of genomic medicine. Likewise, transformation of human medicine by harnessing the gut microbes requires the cataloging of reference microbial genomes and their encoded functional elements. Conventional approaches for microbial genome assembly require microbial isolation and culture. Indeed, with the development of culturomics technology, the number of culturable gut microbes has increased greatly[3-6]. However, the culturable taxa are biased toward specific clades, and a large portion of the human gut microbiome remains unculturable[7-9]. To address this, culture-independent methods of metagenome assembly from whole-metagenomic shotgun sequencing (WMS) data have been developed.

Recently, three independent studies have consecutively released large collections of prokaryotic genomes, including many based on metagenome assembly[8-10]. The metagenome-assembled genomes (MAGs) from these studies were then combined with the genomic information deposited in public databases to generate integrated catalogs of prokaryotic genomes and proteins in the human gut[11], the Unified Human Gastrointestinal Genome (UHGG) and Unified Human Gastrointestinal Protein (UHGP) catalogs, respectively. The UHGG contains 204,938 non-redundant genomes that represent 4,644 prokaryotic species and the UHGP catalogs approximately 95 million unique proteins.

Despite the latest advances, the current human gut microbiome catalog is incomplete, partially because the metagenome sampling is biased for geographical location and lifestyle. Specifically, the UHGG is strongly biased towards fecal samples collected in China, Denmark, Spain, and the US. In the present study, to account for the under-sampling of certain metagenomes, we assembled genomes from fecal samples collected from Korea, India, and Japan. Since the genome assembly of low-abundance species in most human fecal samples may require a much deeper sequencing than usually employed, we performed ultra-deep WMS (>30 Gbp or >100 million read pairs) of 90 fecal samples collected from Korea. We then collected public WMS data for 110 fecal samples from India and 805 fecal samples from Japan. We consequently assembled 29,082 prokaryotic genomes, and combined them

72    with the UHGG genomes to generate the Human Reference Gut Microbiome (HRGM),

73    which substantially expands the list of representative species, genomes, proteins, and single-

74    nucleotide variants (SNVs) in the human gut microbiome. The HRGM is a freely available

75    resource and will be invaluable to therapeutic targeting of the gut microbiota.

76

## Results

78    **Assembly of gut microbial genomes from Korea, India, and Japan**

79    We assembled prokaryotic genomes using an in-house bioinformatics pipeline

80    (**Supplementary Fig. 1a, Methods**), which is more exhaustive than similar approaches[8-11]

81    (**Supplementary Table 1**). For instance, we adopted an ensemble method for binning

82    assembled contigs, as it showed better performance than individual binning tools[12,13]. We

83    hypothesized that metagenomes harbored by individuals from under-represented geographical

84    locations and lifestyles would expand the current catalog of human gut microbiome.

85    Therefore, we performed *de novo* genome assembly of fecal samples from three Asian

86    countries: Korea, India, and Japan (referred to here as KIJ samples, **Supplementary Table 2**).

87    At the start of the current study, WMS data for 805 and 110 fecal samples from Japan and

88    India, respectively, were publicly available but not included in the UHGG[14,15]. To

89    complement these data, we generated WMS data for fecal samples collected from 90 donors

90    recruited in Korea. We set the minimum completeness at 50% and the maximum

91    contamination at 5% for genomes of minimum quality. We divided the genome bins into two

92    groups: high quality (HQ) genomes with ≥90% completeness and ≤5% contamination, and

93    medium quality (MQ) genomes (the remaining genomes). This yielded 29,082 KIJ sample

94    MAGs: 7,767 from Korea, 563 from India, and 20,752 from Japan.

95

96    **Ultra-deep sequencing facilitates the genomic assembly of low-abundance taxa**

97    To investigate the impact of metagenome sequencing depth on *de novo* genome assembly, we

98    performed ultra-deep sequencing of the 90 Korean fecal samples (>30 Gbp or >100 million

99    read pairs); the depth was approximately 5-fold deeper than the normal sequencing depth

100   (**Fig. 1a**). Despite sequencing at the normal depth, fecal samples from Japan had a larger total

101 read length than Korean samples because of a much larger sample size (**Fig. 1b**). For nine of

102 the 90 Korean samples, approximately 60 Gbp was sequenced for the study of sequencing

103 depth effect on genome assembly. We then generated 81 simulated WMS datasets (9 different

104 depths for each of the 9 original samples with ~60 Gbp depth) and used the same pipeline of

105 *de novo* genome assembly for all samples. As expected, the number of HQ and MQ genomes

106 increased with the increasing sequencing depth. However, the growth rate simultaneously

107 decreased and the proportion of HQ genomes became stable after the initial phase of rapid

108 growth (**Fig. 1c**). Next, we investigated whether the increased sequencing depth improved the

109 quality of assembled genomes. We compared the assembly quality of MAGs for the same

110 species in two different simulated samples at adjacent sequencing depths (**Supplementary**

111 **Fig. 2**; **Methods**). The quality of MAGs from the greater sequencing depth was significantly

112 higher than that of genomes from the lower sequencing depth in terms of completeness,

113 contamination, N50, and genome size (**Fig. 1d,e; Supplementary Fig. 3a,b**). However, the

114 degree of improvement of the assembly quality diminished as the sequencing depth increased.

115 We then examined the effect of sequencing depth using the actual WMS data for KIJ samples.

116 The number of HQ and MQ genomes assembled from each sample was highest in the ultra-

117 deep sequenced samples from Korea (**Fig. 1f**). However, the proportion of HQ genomes in

118 samples from Korea and Japan was not significantly different (**Fig. 1g; Supplementary Fig.**

119 **3c**). Notably, the genome assembly yield, i.e., the number of assembled genomes divided by

120 the total sequencing length, was highest for samples from Japan (**Fig. 1g**). This suggests that

121 sequencing hundreds of samples at a depth of 5–10 Gbp may constitute the most effective

122 strategy for cataloging MAGs for a given population.

123 The ultra-deep sequencing may be advantageous for the genome assembly for low-abundance

124 taxa. To test this, we compared MAGs exclusively assembled from each country but not

125 included in the UHGG, i.e., 224, 388, and 18 genomes from Korea, Japan, and India,

126 respectively. We then estimated their relative abundance in fecal samples in an independent

127 population of 926 fecal samples from the US[16], using Kraken2[17]. The genomes assembled

128 exclusively from Korean samples shifted towards low-abundance taxa compared with

129 genomes assembled from samples from other countries (**Fig. 1h**), which confirmed the

130 original hypothesis.

131

132 **Cataloging reference genomes of 5,414 prokaryotic species from the human gut**

133 To construct the most comprehensive reference database for the human gut microbiome, we

134 integrated the newly generated 29,082 MAGs from KIJ samples with the UHGG genomes

135 using dereplication approach (**Supplementary Fig. 1b**, **Methods**). Dereplication of the

136 29,082 MAGs resulted in 2,199 clusters of genomes. We selected a representative genome

137 from each cluster to catalog the genomes for 2,199 representative species, which we then

138 integrated with 4,644 representative genomes from the UHGG, via dereplication, resulting in

139 5,414 clusters of genomes. Finally, we selected 5,414 representative genomes and assigned

140 their phylogenetic classifications using GTDB-Tk[18] (**Fig. 2**). Among these representative

141 genomes, 4,531 (83.7%) genomes were exclusively assembled from metagenomic data,

142 which confirmed the notion that the major portion of the human gut microbiome has not yet

143 been isolated. We identified 16S rRNA sequences in 2,542 representative genomes (47%)

144 (**Supplementary Fig. 4**), covering the majority of phylogenetic clades. Unlike conventional

145 databases of 16S rRNA sequences, the new database provides opportunities for functional

146 interpretation of the detected taxa because it contains genomes corresponding to the 16S

147 rRNA sequences.

148 The inclusion of MAGs from KIJ samples in the new database allowed several improvements

149 on the UHGG. First, we reduced the data bias toward China among Asian countries

150 (**Supplementary Fig. 5a**). Second, we expanded the total number of non-redundant reference

151 genomes by 13.25% and the number of representative species by 16.6% increase

152 (**Supplementary Table 3**). Among the 5,414 representative genomes, 780 genomes were

153 assembled from KIJ samples only, and 536 representative genomes from the UHGG were

154 replaced with new MAGs from KIJ samples. Hence, 1,316 representative genomes (28.3%)

155 were updated in the HRGM (**Supplementary Fig. 5b**).

156

157 **New MAGs from Korea, India, and Japan are associated with diet-related lifestyles**

158 Notably, *Bacteroidaceae* family (**Fig. 3,** redtree branches) was enriched in the updated

159 MAGs ($P < 0.001$, Fisher's exact test). Almost half the genomes from this family are from

160 the *Bacteroides* genus and approximately two-thirds of the other half are from the *Prevotella*

161 genus (**Supplementary Fig. 6**). Interestingly, three widely dispersed regions in the

162    phylogenetic tree were highly enriched in the updated genome set. The first region ("a")

163    encompasses a portion of the *Prevotella* genus and includes 30 genomes annotated as

164    *Prevotella copri*. Accordingly, westernized populations with a typically high-fat and low-

165    complex carbohydrate diet exhibit low prevalence and diversity of *P. copri* compared with

166    non-westernized populations[19]. The second region ("b") encompasses a portion of the

167    *Bacteroides* genus and includes 22 genomes annotated as *Bacteroides plebeius*. This species

168    is typically found in Japanese subjects whose diet includes seaweed-rich food, such as sushi[20].

169    It has been suggested that *B. plebeius* harbors genes encoding an enzyme specific for algal

170    carbohydrates, acquired from marine microbes. The third region ("c") also encompasses a

171    portion of the *Bacteroides* genus and includes 12 genomes annotated as *Bacteroides vulgatus*,

172    which is typically present in the human distal gut, where undigested plant polysaccharides

173    and proteins exist in large quantities[21]. Together, these observations indicate that the new

174    MAGs from KIJ samples are associated with the diet-related lifestyles in Japan and Korea.

175

**SNV density is positively associated with the speciation rate of gut commensals**

177    We then aligned genomes of species clusters containing ≥3 genomes with the representative

178    genome and mapped SNVs (**Methods**). This yielded 274,543,071 SNVs from 2,821 species

179    clusters, representing 10.07% and 13.34% increases, respectively, from the UHGG. The

180    Actinobacteriota phylum had the highest SNV density (**Fig. 3a**). Phylogenetically

181    overdispersed branches of Actinobacteriota species were apparent in both, the HRGM and

182    UHGG. The majority of genomes from the overdispersed tree region belonged to the

183    *Collinsella* genus. We divided these genomes into ones from a tree region with a modest

184    phylogenetic dispersion (MD, 20 genomes) and those with a high phylogenetic dispersion

185    (HD, 619 genomes) (**Fig. 3b**). Although the majority of genomes were not annotated at the

186    species level, *Collinsella aerofaciens* was enriched in the HD group and other known

187    *Collinsella* species were enriched in the MD group (**Fig. 3c**). SNV density in HD group was

188    significantly higher than that of MD group (**Fig. 3d**).

189    SNV, a within-species genetic variation, is a major mechanism for the adaptation of

190    commensal species to a distinct host environment. Wide dispersion of species branches

191    indicates rapid speciation. Accordingly, high SNV density for a species with an overdispersed

192    tree may indicate that the degree of within-species genetic variation may be positively

193    associated with the speciation rate of gut commensals. To test this, we examined the

194    correlation between SNV density of representative species and their phylogenetic distance to

195    the five nearest species. The branch length to the neighboring species in the phylogenetic tree

196    of a species that arose during rapid speciation tends to be short. We observed an inverse

197    correlation between the average phylogenetic distance to the five nearest species and their

198    SNV density (**Fig. 3e**), and a significantly higher SNV density for the top 10% species with

199    shorter phylogenetic distance to the nearest five species than those for the bottom 90%

200    species (**Fig. 3f**). This supports the model of a positive correlation of SNV density and the

201    speciation rate of gut commensals.

202

203    **Functional landscape of 103 million proteins from human gut prokaryotes**

204    Information on proteins encoded in the human gut microbes will facilitate the functional

205    characterization of disease-associated microbiota. Using an in-house computational pipeline

206    for cataloging human gut prokaryotic proteins (**Supplementary Fig. 1c** and **Supplementary**

207    **Fig. 7**), we first identified 64,661,728 CDS (coding sequences) from 29,082 genomes from

208    KIJ samples using Prodigal[22]. To reduce redundancy in the protein catalog, we first executed

209    CD-HIT[23] at 100% similarity level and then combined with proteins cataloged by the UHGP-

210    100[11]. The consolidated protein catalog was next consecutively clustered by CD-HIT at lower

211    sequence similarity levels: 95%, 90%, 70%, and 50%. This led to approximately 103.7, 20.0,

212    14.8, 8.5, and 4.7 million proteins at the sequence similarity levels of 100%, 95%, 90%, 70%,

213    and 50%, respectively.

214    Unexpectedly, we observed that the UHGP contains proteins that are 100% identical, even in

215    a catalog at 50% sequence similarity level. For instance, among the UHGP-50 proteins,

216    GUT_GENOME232012_01109 and GUT_GENOME231777_00918 have an identical amino

217    acid sequence. We identified 8,663, 82,507, 243,362, and 75,620,150 proteins that are

218    redundant at 100% similarity in the UHGP-50, UHGP-90, UHGP-95, and UHGP-100,

219    respectively. Exclusion of the UHGP proteins that were 100% identical revealed that the

220    HRGM contains more proteins than UHGP at all levels of sequence similarity except for 50%

221    (**Supplementary Table 3**).

8

222    To facilitate the functional interpretation of gut microbiome profiles, we next annotated

223    functional genomic elements and proteins in the HRMG. We predicted and annotated non-

224    coding RNAs and functional peptides, using Prokka[24]; antibiotic resistance genes, using

225    RGI[25]; biosynthetic gene clusters, using antiSMASH[26]; and 16S rRNA regions, using

226    barrnap[27]. For functional annotation of proteins, we used eggNOG-mapper[28]. Notably, the

227    landscape of antibiotic resistance ontology revealed that phylogenetically close species in the

228    human gut tend to share antibiotic resistance mechanisms (**Supplementary Fig. 8**). A

229    significantly large portion of the human gut prokaryotic proteins has not yet been functionally

230    annotated. For the HRGM protein catalogs at 100%, 95%, 90%, 70%, and 50% similarity

231    levels, 13.13%, 28.05%, 29.17%, 36.35%, and 47.62% of proteins, respectively, had no

232    functional annotation, according to eggNOG-mapper. This effect appears to be amplified by

233    redundant proteins, resulting in a reduced annotation rate at low similarity level. Further, the

234    annotation rate of proteins that are shared by many species is higher than that of species-

235    specific proteins (**Supplementary Fig. 9**).

236

237    **HRGM improves taxonomic and functional classification of sequencing reads**

238    According to a recent benchmark study, whole-DNA–based methods outperform marker-

239    based methods for taxonomic classification of metagenomic sequencing reads[29]. The

240    performance of whole-DNA–based methods relies on the quality of the reference genome

241    database. The standard databases lack numerous genomes of species that exist in the human

242    gut, which leads to false-negatives, while including many genomes from other microbial

243    communities, which leads to false-positives[29]. We hypothesized that the HRGM, which is

244    specific to the human gut microbiome and more comprehensive than other databases, can

245    improve the taxonomic classification of sequencing reads. We used Kraken2[17] to compare the

246    taxonomic classification of three genome databases: a standard database that contains

247    RefSeq[30] complete genomes (RefSeq CG) of bacterial, archaeal, and viral domains; the

248    UHGG-based database; and the HRGM-based database. To generate independent test

249    datasets, we compiled WMS data for 1,022 fecal samples from the US, Cameroon,

250    Luxembourg, and Korea, which were not included in the UHGG nor HRGM. We then

251    evaluated the efficacy of Kraken2 classification based on the proportion of classified reads

252    (**Methods**). The classification efficacy using the UHGG and HRGM-based databases was

253  substantially higher than that of the standard database (**Fig. 4a,b**, $P < 0.001$, two-sided

254  Wilcoxon signed-rank test). In addition, the variance of the read classification rate of custom

255  databases was significantly smaller than that of the standard database, except for the

256  Cameroon population (**Fig. 4a**, $P < 0.001$, Brown-Forsythe test). Importantly, the

257  classification efficacy of the HRGM-based database was significantly improved compared

258  with that of the UHGG-based database for the four test samples (**Fig. 4a,c**, $P < 0.001$, two-

259  sided Wilcoxon signed-rank test), which suggests that the updated reference genome database

260  improves taxonomic classification of the gut metagenomic sequencing data.

261  Next, we investigated the efficacy of functional classification based on the number of aligned

262  sequencing reads from reference protein databases. Because of the extremely large number of

263  reference proteins, we used only 40 samples randomly selected from the 1,022 fecal samples

264  (10 samples from each population), and aligned the sequencing reads with the UHGP-95 and

265  HRGM-95 protein catalogs (**Methods**). The number of aligned reads was 1.31% higher, on

266  average, with HRGM-95 in all tested samples than with UHGP-95 (**Fig. 4d**), although

267  HRGM-95 contains 0.4% more proteins than UHGP-95.

268  Taken together, the newly assembled genomes from under-represented Asian countries

269  significantly improve the genome and protein databases for metagenomic analysis of both,

270  taxonomic and functional profiling.

271

**Reliable taxonomic profiling of low-abundance taxa requires deep sequencing**

273  Taxonomic profiles obtained by shallow sequencing (0.5–2 million reads) highly correlate

274  with those obtained by ultra-deep sequencing (2.5 billion reads)[31]. However, this evaluation

275  is based on entire taxa, in which highly abundant or core taxa govern the correlation measure.

276  Further, low-abundance taxa likely play important, as yet unknown, biological roles in the gut

277  microbial communities[32,33]. We therefore evaluated the impact of sequencing depth on the

278  reliability of taxonomic profiling for different ranges of taxon abundance. We generated a

279  simulated dataset at various sequencing depths 16 new Korean fecal samples, and not

280  included in the HRGM. We then stratified the taxonomic features into eight different groups,

281  according to the mean relative abundance (**Fig. 5a,b**). We calculated the mean Pearson

282  correlation coefficient (*PCC*) and the mean Spearman correlation coefficient (*SCC*) between

283    the taxonomic profiles at different sequencing depths for different mean relative abundances

284    (**Methods**). The taxonomic profile similarity between two groups showed increasing *PCC*

285    and *SCC* with an increasing sequencing depth. For example, >10 million read pairs (3 Gbp)

286    may need to have taxonomic profiles that highly correlate (*PCC* > 0.9) with those based on

287    80 million read pairs (25 Gbp) to account for the features with lowest 13.92% of relative

288    abundance (relative abundance < 1e–06) (**Fig. 5c and Supplementary Fig. 10a**). For *SCC* >

289    0.9, the required sequencing depth increased to 20 million read pairs (6 Gbp) for taxonomic

290    features with a similar level of relative abundance (**Fig. 5b and Supplementary Fig. 10b**).

291    Overall, these observations suggest that deep sequencing (>20 million read pairs) may be

292    required to obtain reliable taxonomic profiles of low-abundance taxa.

293

294    **Sequencing 30 Gbp is optimal for functional profiling of the human gut microbiome**

295    Next, using the protein catalog, we investigated the optimal sequencing depth for functional

296    profiling of the human gut microbiome. Since the detection of gene content generally requires

297    a much deeper sequencing depth than that for the detection of genomes, we analyzed the

298    WMS data for five Korean fecal samples at a depth of approximately 200 million read pairs

299    (60 Gbp) (**Methods**). The number of the detected coding genes initially grew rapidly as the

300    sequencing depth increased, but later approached the estimated maximum count

301    (**Supplementary Fig. 11a**). The curves fitted well ($R^2$ > 0.99) two-site saturation models[34],

302    and we hence estimated the maximum number of coding genes for each sample using the

303    regression model. Interestingly, the estimated maximum gene counts in the samples differed,

304    reflecting the different alpha diversity of the microbial community. However, all samples

305    showed very similar normalized maximum gene count curves, with over 80% of the gut

306    microbial coding genes detected by sequencing 30 Gbp or 100 million read pairs in all

307    samples (**Supplementary Fig. 11b**). Sequencing another 30 Gbp would fail to detect 90% of

308    the maximum gene count. Therefore, 100 million read pairs is the optimal sequencing depth

309    for the best trade-off between the sequencing cost and the gain-of-functional information for

310    WMS-based studies of the human gut microbiome.

311

312    **Profiling cross-reactivity potential identifies autoimmune-associated commensals**

313   Microbial peptides homologous to the host auto-antigens may stimulate host immune cells

314   and, hence, the hypothesis of molecular mimicry has emerged as a mechanism underlying

315   autoimmune diseases[35]. To systematically evaluate this hypothesis, we mapped microbial

316   peptide sequences homologous to the human self-antigens involved in autoimmune diseases

317   onto the genomes of HRGM representative species. We first compiled autoimmune disease-

318   related antigen set from the Immune Epitope Database (IEDB)[36], and then used it for

319   homology-searches of microbial peptide sequences from 5,414 representative species. We

320   thus identified species with a high cross-reactivity potential based on the density of the

321   encoded cross-reactive epitopes. Because the number of epitope-containing genes (ECG)

322   increased as the number of coding genes increased (**Fig. 6a**), we divided the ECG count by

323   the total number of genes for each species. Some human gut commensals had a relatively

324   high cross-reactivity potential (**Fig. 6b,c, Methods**). On the genus level, *Akkermansia*,

325   *Alistipes*, *Bifidobacterium*, *Lawsonibacter*, *Oscillibacter*, *Prevotella*, and *Sutterella* have a

326   high cross-reactivity potential (**Fig. 6d**). Indeed, many of them are associated with

327   autoimmune diseases. For example, *Akkermansia muciniphila* is abundant in the enthesitis-

328   related arthritis patients[37], while *Bifidobacterium* is enriched in these[37] and inflammatory

329   bowel disease (IBD) patients[38]. Increased abundance of *Oscillibacter* is accompanied by

330   increased levels of interleukin 6[39], a pro-inflammatory cytokine that can disrupt the immune

331   homeostasis and increase the risk of autoimmune diseases. The abundance of intestinal

332   *Prevotella copri* is strongly correlated with the risk of arthritis[40] and *Sutterella*

333   *wadsworthensis* is enriched in ulcerative colitis patients who do not respond to fecal

334   microbiota transplantation[41]. These suggests that cross-reactivity potential of commensal

335   genomes is predictive for human gut microbiota associated with autoimmune diseases.

336

## Discussion

338   In the present study, we constructed an improved catalog of the human reference gut

339   prokaryotic genomes and their proteins, by including MAGs from fecal metagenomes from

340   under-represented Asian countries. Inclusion of the newly assembled genomes expanded the

341   catalog size by over 10%. In addition, we demonstrated that database expansion also

342   significantly improved the taxonomic and functional classification of sequencing reads. Many

343   new MAGs were associated with diet-related lifestyles at the sampled geographic locations.

344    Therefore, complementation of metagenome datasets to account for under-sampled
345    geographical locations and lifestyles might be an effective strategy for improving the human
346    reference gut microbiome.

347    We also demonstrated that the analysis of microbial DNA and peptide sequences facilitates
348    the understanding of gut commensal speciation and interactions with the host immunity. The
349    colonizing commensal microbes adjust to their host environment via genetic changes and
350    selection, which lead to genetic variation within species. We cataloged the SNVs of
351    conspecific genomes and found that the SNV density of gut prokaryotic species is inversely
352    correlated with the phylogenetic distance to their neighboring species. This may suggest that
353    the degree of within-species genetic variation is positively associated with the speciation rate
354    of gut commensal microbes. Whether SNV actually enhances the speciation rate should be
355    addressed in future investigations. Finally, we showed that systematic analysis of microbial
356    peptide sequences homologous to the host auto-antigens allows the prediction of gut
357    microbial taxa potentially associated with autoimmune disease via the mechanism of
358    molecular mimicry. Such analysis is only possible if microbial protein sequences are
359    available with the corresponding taxonomic information.

360    As the WMS analysis for population-wide human gut microbiome profiling increases in
361    popularity, the choice of sequencing depth is an important factor to consider in study design.
362    Here, we demonstrated that deep sequencing (>20 million read pairs) is necessary for reliable
363    taxonomic profiling of low-abundance commensals. The current knowledge of human gut
364    microbiome is biased towards core taxa that are usually highly abundant. Low sequencing
365    depth (e.g., 0.5–2 million read pairs) may be sufficient for the profiling of core taxa, but not
366    those with low abundance. Deep sequencing may therefore be required for the WMS-based
367    analysis of human gut microbiome to investigate the function of relatively unexplored low-
368    abundance species. Accordingly, the current study provides the guidelines for the choice of
369    sequencing-depth for the analysis of human gut microbiome for different purposes.

370    In conclusion, the HRGM database, which contains information on various biological entities,
371    from DNA and protein sequences to pan-genomes of species, is a versatile resource for
372    functional dissection of disease-associated gut microbiota. The data will be available via a
373    web server (www.mbiomenet.org/HRGM/) and will be periodically updated as new WMS
374    data for fecal samples become publically available.

375

## Methods

### Sequencing fecal metagenome samples from Korea, India, and Japan

376

377

378 WMS data for fecal samples from India and Japan were obtained from published studies[14,15].

379 Fecal WMS data for India were generated for 110 healthy donors in North-Central and

380 Southern India[14]. Although the sequencing depth was relatively low (1.2 Gbp on average), it

381 was expected that many novel genomes would be assembled because MAGs from India are

382 not included in the existing catalogs. By contrast, 805 MAGs from Japan are included in the

383 UHGG. However, it was expected that the inclusion of the recently published deep-

384 sequencing WMS data for 645 Japanese fecal samples (6.5 Gbp on average)[15] would greatly

385 expand the number of MAGs for Japan. In addition, ultra-deep WMS data (31 Gbp on

386 average) were generated for fecal samples from 90 Koreans recruited by the Severance

387 Hospital (Seoul, Korea; IRB No 4-2020-0309 and IRB No 4-2017-0788). Written informed

388 consent was obtained before the study. The UHGG does not contain any MAGs from Korea.

389 The libraries were prepared as described in the TruSeq Nano DNA Library Prep Reference

390 Guide (Illumina #15041110). Briefly, 100 ng DNA was fragmented using LE220 Focused

391 ultrasonicator (Covaris, Inc.). Fragmented DNA was end-repaired and approximately 350-bp

392 fragments were obtained after size selection. After adapter ligation, eight PCR cycles were

393 performed. Library quantification was performed as described in the Kapa Illumina Library

394 Quantification Kit (Kapa Biosystems, #KK4854). Next, 150 bp ×2 paired-end sequencing

395 was performed using Illumina HiSeq4000. In summary, new WMS data for 845 fecal

396 samples collected from Korea, India, and Japan were obtained. The total read length was 7.2

397 Tbp. All samples used in the current study are described in **Supplementary Table 2**.

398

### Metagenome assembly and binning

399

400 The adapter sequences were trimmed, and low-quality bases and short reads were removed

401 from WMS data using Trimmomatic v0.39[42]. Next, the reads were aligned with the human

402 genome GRCh38.p7 using Bowtie2 v2.3.5[43], and the aligned reads were then removed. The

403 majority of quality-controlled reads were assembled as contigs using metaSPAdes[44], which is

404     a metagenome-specific pipeline of SPAdes v3.13.0. For unknown reasons, and regardless of

405     sample size, metaSPAdes runtime was excessively long for 107 samples. In those cases,

406     MEGAHIT v1.2.8[45] was used (**Supplementary Table 2**).

407     Genome bins were generated using the ensemble approach and three binning tools:

408     MetaBAT2 v2.13[46], MaxBin2.0 v2.2.6[47], and CONCOCT v1.1.0[48]. First, the reads from each

409     sample were first aligned with the assembled contigs from the previous step using Bowtie2,

410     and the three binning programs were initiated. The minimum size of a contig for binning was

411     set at 1,000 bp, except for MetaBAT2, which requires at least 1,500 bp. The three binning

412     predictions were combined for improved binning results using the bin refinement module of

413     MetaWRAP v1.2.2[12], which uses CheckM v1.0.18[49] to evaluate the quality of genome bins in

414     terms of completeness and contamination rate. The minimum completeness was set at 50%,

415     the maximum contamination at 5%, and the minimum quality score ($Completeness - 5 \times$

416     $Contamination$) at 50. The same threshold values for CheckM results were applied during the

417     construction of the UHGG. This resulted in 7,767 genomes from Korean samples, 563

418     genomes from Indian samples, and 20,752 genomes from Japanese samples (29,082 genomes

419     in total). The genome bins were divided into two groups: HQ, bins with over 90%

420     completeness and less than 5% contamination; and MQ, bins with 50–90% completeness and

421     less than 5% contamination.

422

423     **Generation of genomic species clusters**

424     Groups of genomes that corresponded to species were generated using a two-step iterative

425     procedure. Preliminary clustering was performed using Mash v2.2[50] algorithm. Mash

426     distances were calculated for all possible pairs of genomes using the "-s 10,000" parameter.

427     Next, the average-linkage–based hierarchical clustering was performed, at a cutoff of 0.2.

428     Mash algorithm is sufficiently fast to calculate all-by-all distances for hundreds of thousands

429     of genomes in a timely manner. However, this compromises the accuracy, especially for low-

430     coverage genome pairs[51], which are common in MAGs. Therefore, to improve cluster quality,

431     ANImf[51] was calculated for every pair of genomes within each initial cluster. To avoid the

432     over-estimation of ANI by local alignment, a minimum coverage threshold was applied for

433     each pair. The coverage cutoff of genome A and genome B was determined at *min(0.8,*

434  *Completeness of genome A × Completeness of genome B)*. If the alignment coverage between

435  two genomes was lower than the cutoff, they were regarded as different genomes. The

436  genomes were then clustered using the average linkage-based hierarchical clustering at a

437  cutoff of 0.05 (or 95% identity), which is a widely accepted ANI threshold for species-level

438  boundary[4,9-11,52]. The genome intactness score $(S)$[9,11], $S = Completeness - 5 \times$

439  *Contamination* $+ 0.5 \times log_{10}(N50)$, was then calculated. For clusters containing more than

440  two genomes, a genome with the highest $S$ was selected as the representative genome for the

441  cluster. The above two-step procedure was iterated until the clusters ceased to change. Hence,

442  2,199 species clusters were generated for 29,082 genomes from KIJ samples, with eight

443  iterations of the aforementioned procedure (**Supplementary Fig. 1a**). Finally, the 2,199

444  genomes were combined with 4,644 genomes from the UHGG, generating 5,414 species

445  clusters for the HRGM at the fourth iteration (**Supplementary Fig. 1b**).

446

447  **Non-redundant genome counting**

448  To count the number of non-redundant genomes, the redundant genomes were removed,

449  similar to what was done for the UHGG pipeline[11]. First, the pairwise genome distance was

450  calculated using Mash[50] and the entire genomes were clustered using average-linkage–based

451  hierarchical clustering, with a 0.001 cutoff (Mash ANI 99.9%). To reduce the computation

452  time, the hierarchical clustering was performed only for the connected components with the

453  distance of 0.1, because it is highly unlikely that genomes that are not within the distance of

454  0.1 are clustered together by a distance of 0.001. In the process, 22,761 genomes were

455  clustered into 8,508 conspecific genome bins. Multiple genomes from the same sample in the

456  same species bin were counted only once.

457

458  **Taxonomic and functional annotation of representative species genomes**

459  The taxonomic annotation of 5,414 representative species genomes was performed using the

460  "classify_wf" function of GTDB-Tk v1.0.2[18]. The reference version was GTDB R04-RS89,

461  released in June 2019. Genomic features, such as CDS, rRNA, and tRNA, were identified and

462  annotated in each genome using Prokka v1.14.5[24] with "--kingdom Bacteria" and "--kingdom

463  Archaea" parameters for the bacterial and archaeal genomes, respectively. With the protein

464    sequences predicted by Prokka, the antibiotic resistance genes were annotated using RGI

465    v5.1.0[25] with default parameters. The landscape of antibiotic resistance potential of 5,414

466    species-representative genomes is depicted in **Supplementary Fig. 8**. Finally, the secondary

467    metabolite gene cluster was annotated using antiSMASH v5.1.2[26]. For the full-featured

468    annotation, the "--cb-general, --cb-knownclusters, --cb-subclusters, --asf, --pfam2go, --

469    smcog-trees, --cf-create-clusters" parameters were set.

470    To render the HRGM useful for the 16S rRNA sequencing-based metagenomic analysis, the

471    16S rRNA regions for 5,414 representative species genomes were predicted using barrnap

472    v0.9[27] tool and the "--evalue 1e-05" parameter, and "--kingdom bac" and "--kingdom arc"

473    parameters for bacterial and archaeal genomes, respectively. The 16S rRNA sequences were

474    thus directly predicted from 1,364 representative species genomes. For the remaining 4,050

475    representative species, the search for 16S rRNA sequences was expanded to their conspecific

476    genomes. The barrnap analysis was used for the genomes from KIJ samples and pre-

477    established 16S rRNA region annotations were used for the genomes from the UHGG.

478    Within the expanded search space, 16S rRNA sequences were identified for 1,178 additional

479    genomes. Consequently, 16S rRNA sequences were generated for 2,542 species in the

480    HRGM (**Supplementary Fig. 4**).

481

482    **Cataloging SNVs**

483    For the species bins with more than three genomes, SNVs were identified using the codes

484    provided by the UHGG[11]. Briefly, non-representative genomes were aligned with the

485    representative genome in the species bin using nucmer 4.0.0beta2[53]. Best bi-directional

486    alignments were identified using the delta-filter program and "-q –r" options, and SNVs were

487    annotated using the show-snp program; nucmer, delta-filter, and show-snp are software

488    packages of MUMmer v3[54]. For each species bin ($G$) whose representative genome is $r$, the

489    number of SNV per kb was calculated as follows:

490
$$SNV\ per\ kb = \frac{\sum_{g \in (G-\{r\})} \frac{\#SNV_{r,g}}{Aligned\ length_{r,g}/1000}}{n(G)-1}$$

491    *SNV per kb* was only calculated for 1,521 species bins with ≥10 genomes to avoid bias. For

492    the 1,521 genomes, the average phylogenetic distance to the five nearest species was

493     calculated using the IQ-Tree[55].

494

**495     Cataloging gut prokaryotic proteins and their functional annotation**

496     Overall, 64,661,728 CDS were identified in 29,082 genomes from the KIJ set using Prodigal

497     v2.6.3[22] and "-c -m -p single" parameters. Since many proteins were derived from conspecific

498     genomes, the catalog may have included many homologous proteins. To reduce the

499     redundancy in the protein catalog, CD-HIT v4.8.1[23] was adopted. To reduce CD-HIT running

500     time, identical proteins were first clustered and then CD-HIT was executed at 100%

501     similarity level. The cataloged proteins were then combined with those in UHGP-100[11]. The

502     consolidated protein catalog was subsequently submitted to CD-HIT clustering analysis at

503     five different sequence similarity levels, 100%, 95%, 90%, 70%, and 50%. For accurate and

504     efficient clustering, a multi-step iterative clustering method recommended by the CD-HIT

505     tutorial was adopted. For instance, the CD-HIT-95 protein catalog (a 95% similarity level

506     protein catalog) was constructed based on CD-HIT-100 proteins, and the CD-HIT-90 protein

507     catalog was constructed based on CD-HIT-95 proteins. This resulted in approximately 103.7

508     million, 20.0 million, 14.8 million, 8.5 million, and 4.7 million proteins at the sequence

509     similarity levels of 100%, 95%, 90%, 70%, and 50%, respectively. The overall pipeline for

510     protein catalog construction is depicted in **Supplementary Fig. 7**.

511     Representative protein sequences in the five protein catalogs were functionally annotated

512     using eggNOG-mapper v2.0.1[28], which is based on the eggNOG protein database v5.0[56]. The

513     resultant annotations include eggNOG orthologs and functional terms from several databases,

514     including Gene Ontology (GO)[57] and Kyoto Encyclopedia of Genes and Genomes (KEGG)[58].

515     Further, for each protein cluster, taxonomic origins of all member proteins and the lowest

516     common ancestor of the cluster were tracked and annotated.

517     The numbers of shared species and shared phyla of proteins in the HRGM-50 protein catalogs

518     were annotated based on the taxonomic annotation of member proteins. The number of

519     shared species was binned at the bin size of 10, then the annotation rate for each protein bin

520     was calculated as the number of annotated proteins divided by the number of proteins in the

521     bin.

522

523 **Reconstruction of the phylogenetic tree**

524 For the bacterial and archaeal genomes, 120 and 122 universal marker genes, respectively,

525 were predicted by the GTDB-Tk[18]. Using the concatenated sequences of marker genes, the

526 maximum-likelihood tree was generated using IQ-TREE[55]. The phylogenetic tree of bacterial

527 genomes was visualized using iTOL[59].

528

529 **Kraken2 databases**

530 The Kraken2 v2.0.8-beta[17] custom database for the HRGM representative genomes was

531 prepared based on the taxonomic annotations in GTDB-TK[18]. When two or more genomes

532 were annotated to the same taxon, they were discriminated at the succeeding lower rank. For

533 example, if *genome a* and *genome b* were both annotated to *species_A*, *genome a* and *genome*

534 *b* were annotated as *Species_A;strain_1* and *Species_A;strain_2*, respectively. By doing so,

535 the user can select a taxonomic rank, thereby measuring species abundances together or

536 individually.

537 The Kraken2 database for the UHGG[11] was downloaded from UHGG FTP on March 6, 2020.

538 The Kraken2 standard database was downloaded and constructed using "kraken2-build --

539 standard" command on July 14, 2020.

540

541 **Measuring taxonomic classification rate of sequencing reads**

542 WMS data were compiled for publicly available data for 926, 54, and 26 fecal samples from

543 the US[16], Cameroon[60], and Luxembourg[61,62], respectively. WMS data for 16 fecal samples

544 collected from Korea, which were not included in the HRGM, were also used. These 1,022

545 fecal samples were neither used for the UHGG nor for the HRGM. The data were pre-

546 processed and taxonomically classified using Kraken2 with standard database, UHGG-based

547 database, and HRGM-based database. The taxonomic classification rate was then calculated

548 based on the proportion of aligned sequence reads in a sample with respect to the database.

549

550 **Measurement of functional classification rate of sequencing reads**

551   The functional classification rate of sequencing reads was determined based on the number of

552   aligned reads against the protein catalog. For the analysis, WMS data were randomly selected

553   for ten fecal samples from the Cameroon, Korea, US, and Luxembourg cohorts (the same

554   samples were used for the taxonomic classification assessment). After pre-processing, 40

555   samples were aligned with the UHGP-95 and HRGM-95 protein databases using blastx of

556   DIAMOND v0.9.35.136 [63]. The results were filtered at >80% query coverage (read coverage)

557   and >95% alignment identity thresholds. A pair of reads was treated as two independent reads.

558   For multiple alignments of a read, only the best alignments by bit score and e-value were

559   considered.

560

561   **Finding the optimal sequencing depth for gene-level analysis of the gut microbiome**

562   For five Korean fecal samples, WMS data generated at a sequencing depth of >60 Gbp, the

563   reads were aligned against the HRGM-95 protein database using blastx of DIAMOND[63].

564   Alignment results with >80% read coverage and 80% identity were included in further

565   analysis. For each sample, the number of detected genes with at least one aligned read was

566   counted by iteratively removing 1000 randomly selected reads. The number of the detected

567   genes for a given sequencing depth exhibited a saturation curve. The curve fitted well ($R^2$ >

568   0.99 for all samples) the two-site binding model[34]. The required sequencing depth for a given

569   gene coverage was determined based on the estimated maximum number of genes according

570   to the equation.

571

572   **Evaluation of the effect of sequencing depth on *de novo* genome assembly**

573   Nine Korean samples with sequencing depth of >52.5 Gbp (**Supplementary Table 2**) were

574   selected for analysis. Then, 0.5, 2.5, 5, 10, 20, 40, 80, 125, and 175 million read pairs were

575   randomly sampled from each of these samples. As the average read-pair length was 300 bp,

576   the sequencing depths of these random samples corresponded to 150 Mbp, 750 Mbp, 1.5 Gbp,

577   3 Gbp, 6 Gbp, 12 Gbp, 24 Gbp, 37.5 Gbp, and 52.5 Gbp, respectively (**Supplementary Fig.**

578   **2**). For the 81 simulated samples (9 samples × 9 depths), *de novo* genome assembly was

579   performed using the same pipeline as that used for the database construction.

580   Two adjacent sequencing depths (e.g., 125 *vs*. 175 million read pairs) were compared to

581 evaluate the effect of sequencing depth on the *de novo* genome assembly. Samples with a
582 greater sequencing depth may yield more MAGs with over 50% completeness, yet with a
583 lower average quality, than those with a lower sequencing depth because of MAGs that
584 barely pass the completeness threshold. Therefore, instead of the average quality scores of all
585 assembled genomes, two genomes assembled at different sequencing depths for the same
586 species clusters were compared. Mash[50] clustering of genomes from two random samples was
587 performed for a comparison based on the average-linkage–based hierarchical clustering, at a
588 threshold of 0.1 (90% identity). Mash clustering was sufficient for clustering conspecific
589 genomes in the simulated samples. Indeed, no cluster had more than two genomes from the
590 same sequencing depth. The assembly quality (completeness, contamination, N50, and
591 genome size) of conspecific genomes at adjacent sequencing depths was then compared.

592

**Evaluation of the effect of sequencing depth on taxonomic profiling**

594 To avoid overestimation of performance, WMS data for 16 Korean fecal samples that have
595 not been used for the HRGM construction and generated at a sequencing depth of >24.5 Gbp
596 were used. From each of the 16 samples, 1, 5, 10, 20, 40, 60, and 80 million read pairs that
597 corresponded to 300 Mbp, 1.5 Gbp, 3 Gbp, 6 Gbp, 12 Gbp, 18 Gbp and 24 Gbp, respectively,
598 were randomly sampled. Taxonomic profiling was then conducted using Kraken2 and the
599 HRGM-based database. Based on the hypothesis that profiling of low-abundance taxa is more
600 affected by sequencing depth than abundant ones, the taxonomic features were stratified at
601 eight different levels of relative abundance, ranging from 1e–07 to 1 with every ten-fold
602 increase (**Fig. 5a,b**). *PCC* and *SCC* between the taxonomic profiles at different sequencing
603 depths were then calculated for each group of features for different levels of relative
604 abundance.

605

**Profiling cross-reactivity potential of the gut prokaryotic genomes**

607 Epitope sequences from autoimmune disease-related self-antigen were compiled from
608 IEDB[36]. "Epitope: Linear epitope", "Antigen: Organism: Homo sapiens", "Host: Homo
609 sapiens", and "Disease: Autoimmune Disease" filters from the IEDB web portal were applied.
610 Epitope sequences that required post-translational modification (e.g., citrullination and

611  deamination) and epitopes shorter than five amino acids were removed. Next, 24,461 unique

612  epitope sequences were aligned with the protein sequences encoded by 5,414 species

613  representatives using BLASTP [64]. For meticulous alignment of short peptide sequences, "-

614  word_size 4", "-evalue 10000", and "-max_target_seqs 100000" options were applied. For

615  every epitope-to-gene pairwise alignment, the Alignment Score (*AS*) was calculated, as

616  follows:

617  $$AS = (\text{match length - gap length}) / \text{epitope length}$$

618  *AS* = 1 alignments were used and the number of protein-coding genes of autoimmune disease

619  epitopes was calculated for every representative species. The number of ECGs was positively

620  correlated with the number of genes. Therefore, the number of ECGs was normalized to the

621  number of genes. To identify epitope-enriched taxonomic clades, EGC per gene of each

622  taxonomic group were compared with the entire 5,414 genomes, and Mann–Whitney P-

623  values and fold-change were calculated.

624

## Data availability

626  Raw metagenomic sequencing data are available from the Sequence Read Archive (accession

627  number will be released upon publication). By accessing the web server,

628  www.mbiomenet.org/HRGM/, users can browse and download all genomes for 5,414

629  representative species, their annotations, and metadata, including geographical origin,

630  taxonomy, genomic content, and genome statistics. The five classes of protein catalogs, 16S

631  rRNA sequences, and SNVs are also provided with their functional annotation and taxonomic

632  origin.

633

## References

635  1  Shreiner, A. B., Kao, J. Y. & Young, V. B. The gut microbiome in health and in
636     disease. *Curr Opin Gastroenterol* **31**, 69-75, doi:10.1097/MOG.0000000000000139
637     (2015).

638  2  Thursby, E. & Juge, N. Introduction to the human gut microbiota. *Biochem J* **474**,
639     1823-1836, doi:10.1042/BCJ20160510 (2017).

640  3  Zou, Y. *et al.* 1,520 reference genomes from cultivated human gut bacteria enable

641    functional microbiome analyses. *Nat Biotechnol* **37**, 179-185, doi:10.1038/s41587-
642    018-0008-8 (2019).

643  4   Forster, S. C. *et al.* A human gut bacterial genome and culture collection for improved
644       metagenomic analyses. *Nat Biotechnol* **37**, 186-192, doi:10.1038/s41587-018-0009-7
645       (2019).

646  5   Poyet, M. *et al.* A library of human gut bacterial isolates paired with longitudinal
647       multiomics data enables mechanistic microbiome research. *Nat Med* **25**, 1442-1452,
648       doi:10.1038/s41591-019-0559-3 (2019).

649  6   Browne, H. P. *et al.* Culturing of 'unculturable' human microbiota reveals novel taxa
650       and extensive sporulation. *Nature* **533**, 543-546, doi:10.1038/nature17645 (2016).

651  7   Eckburg, P. B. *et al.* Diversity of the human intestinal microbial flora. *Science* **308**,
652       1635-1638, doi:10.1126/science.1110591 (2005).

653  8   Almeida, A. *et al.* A new genomic blueprint of the human gut microbiota. *Nature* **568**,
654       499-504, doi:10.1038/s41586-019-0965-1 (2019).

655  9   Nayfach, S., Shi, Z. J., Seshadri, R., Pollard, K. S. & Kyrpides, N. C. New insights
656       from uncultivated genomes of the global human gut microbiome. *Nature* **568**, 505-
657       510, doi:10.1038/s41586-019-1058-x (2019).

658  10  Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by
659       Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle.
660       *Cell* **176**, 649-662 e620, doi:10.1016/j.cell.2019.01.001 (2019).

661  11  Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human
662       gut microbiome. *Nat Biotechnol*, doi:10.1038/s41587-020-0603-3 (2020).

663  12  Uritskiy, G. V., DiRuggiero, J. & Taylor, J. MetaWRAP-a flexible pipeline for
664       genome-resolved    metagenomic    data    analysis.    *Microbiome*    **6**,    158,
665       doi:10.1186/s40168-018-0541-1 (2018).

666  13  Sieber, C. M. K. *et al.* Recovery of genomes from metagenomes via a dereplication,
667       aggregation and scoring strategy. *Nat Microbiol* **3**, 836-843, doi:10.1038/s41564-018-
668       0171-1 (2018).

669  14  Dhakan, D. B. *et al.* The unique composition of Indian gut microbiome, gene
670       catalogue, and associated fecal metabolome deciphered using multi-omics approaches.
671       *Gigascience* **8**, doi:10.1093/gigascience/giz004 (2019).

672  15  Yachida, S. *et al.* Metagenomic and metabolomic analyses reveal distinct stage-
673       specific phenotypes of the gut microbiota in colorectal cancer. *Nat Med* **25**, 968-976,
674       doi:10.1038/s41591-019-0458-7 (2019).

675  16  Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory
676       bowel diseases. *Nature* **569**, 655-662, doi:10.1038/s41586-019-1237-9 (2019).

677    17    Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2.
678        *Genome Biol* **20**, 257, doi:10.1186/s13059-019-1891-0 (2019).

679    18    Chaumeil, P. A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to
680        classify genomes with the Genome Taxonomy Database. *Bioinformatics*,
681        doi:10.1093/bioinformatics/btz848 (2019).

682    19    Thomas, A. M. *et al.* Metagenomic analysis of colorectal cancer datasets identifies
683        cross-cohort microbial diagnostic signatures and a link with choline degradation. *Nat*
684        *Med* **25**, 667-678, doi:10.1038/s41591-019-0405-7 (2019).

685    20    Ledford, H. A genetic gift for sushi eaters. *Nature*, doi:10.1038/news.2010.169 (2010).

686    21    Xu, J. *et al.* Evolution of symbiotic bacteria in the distal human intestine. *PLoS Biol* **5**,
687        e156, doi:10.1371/journal.pbio.0050156 (2007).

688    22    Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site
689        identification. *BMC Bioinformatics* **11**, 119, doi:10.1186/1471-2105-11-119 (2010).

690    23    Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large sets of
691        protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659,
692        doi:10.1093/bioinformatics/btl158 (2006).

693    24    Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-
694        2069, doi:10.1093/bioinformatics/btu153 (2014).

695    25    Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the
696        comprehensive antibiotic resistance database. *Nucleic Acids Res* **48**, D517-D525,
697        doi:10.1093/nar/gkz935 (2020).

698    26    Blin, K. *et al.* antiSMASH 5.0: updates to the secondary metabolite genome mining
699        pipeline. *Nucleic Acids Res* **47**, W81-W87, doi:10.1093/nar/gkz310 (2019).

700    27    Seemann, T. barrnap 0.9: rapid ribosomal RNA prediction.

701    28    Huerta-Cepas, J. *et al.* Fast Genome-Wide Functional Annotation through Orthology
702        Assignment by eggNOG-Mapper. *Mol Biol Evol* **34**, 2115-2122,
703        doi:10.1093/molbev/msx148 (2017).

704    29    Ye, S. H., Siddle, K. J., Park, D. J. & Sabeti, P. C. Benchmarking Metagenomics
705        Tools for Taxonomic Classification. *Cell* **178**, 779-794,
706        doi:10.1016/j.cell.2019.07.010 (2019).

707    30    O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: current status,
708        taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**, D733-745,
709        doi:10.1093/nar/gkv1189 (2016).

710    31    Hillmann, B. *et al.* Evaluating the Information Content of Shallow Shotgun
711        Metagenomics. *mSystems* **3**, doi:10.1128/mSystems.00069-18 (2018).

712  32    Claussen, J. C. *et al.* Boolean analysis reveals systematic interactions among low-
713        abundance species in the human gut microbiome. *PLoS Comput Biol* **13**, e1005361,
714        doi:10.1371/journal.pcbi.1005361 (2017).

715  33    Benjamino, J., Lincoln, S., Srivastava, R. & Graf, J. Low-abundant bacteria drive
716        compositional changes in the gut microbiota after dietary alteration. *Microbiome* **6**, 86,
717        doi:10.1186/s40168-018-0469-5 (2018).

718  34    Wang, Z. X. & Jiang, R. F. A novel two-site binding equation presented in terms of
719        the total ligand concentration. *FEBS Lett* **392**, 245-249, doi:10.1016/0014-
720        5793(96)00818-6 (1996).

721  35    Zhang, X., Chen, B. D., Zhao, L. D. & Li, H. The Gut Microbiota: Emerging
722        Evidence in Autoimmune Diseases. *Trends Mol Med* **26**, 862-873,
723        doi:10.1016/j.molmed.2020.04.001 (2020).

724  36    Vita, R. *et al.* The Immune Epitope Database (IEDB): 2018 update. *Nucleic Acids Res*
725        **47**, D339-D343, doi:10.1093/nar/gky1006 (2019).

726  37    Stoll, M. L. *et al.* Altered microbiota associated with abnormal humoral immune
727        responses to commensal organisms in enthesitis-related arthritis. *Arthritis Res Ther* **16**,
728        486, doi:10.1186/s13075-014-0486-0 (2014).

729  38    Wang, W. *et al.* Increased proportions of Bifidobacterium and the Lactobacillus group
730        and loss of butyrate-producing bacteria in inflammatory bowel disease. *J Clin
731        Microbiol* **52**, 398-406, doi:10.1128/JCM.01500-13 (2014).

732  39    Claesson, M. J. *et al.* Gut microbiota composition correlates with diet and health in
733        the elderly. *Nature* **488**, 178-184, doi:10.1038/nature11319 (2012).

734  40    Scher, J. U. *et al.* Expansion of intestinal Prevotella copri correlates with enhanced
735        susceptibility to arthritis. *Elife* **2**, e01202, doi:10.7554/eLife.01202 (2013).

736  41    Paramsothy, S. *et al.* Specific Bacteria and Metabolites Associated With Response to
737        Fecal Microbiota Transplantation in Patients With Ulcerative Colitis.
738        *Gastroenterology* **156**, 1440-1454 e1442, doi:10.1053/j.gastro.2018.12.001 (2019).

739  42    Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina
740        sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170
741        (2014).

742  43    Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat
743        Methods* **9**, 357-359, doi:10.1038/nmeth.1923 (2012).

744  44    Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new
745        versatile metagenomic assembler. *Genome Res* **27**, 824-834,
746        doi:10.1101/gr.213959.116 (2017).

747  45    Li, D., Liu, C. M., Luo, R., Sadakane, K. & Lam, T. W. MEGAHIT: an ultra-fast

748        single-node solution for large and complex metagenomics assembly via succinct de
749        Bruijn graph. *Bioinformatics* **31**, 1674-1676, doi:10.1093/bioinformatics/btv033
750        (2015).

751  46    Kang, D. D. *et al.* MetaBAT 2: an adaptive binning algorithm for robust and efficient
752        genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359,
753        doi:10.7717/peerj.7359 (2019).

754  47    Wu, Y. W., Simmons, B. A. & Singer, S. W. MaxBin 2.0: an automated binning
755        algorithm to recover genomes from multiple metagenomic datasets. *Bioinformatics* **32**,
756        605-607, doi:10.1093/bioinformatics/btv638 (2016).

757  48    Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition. *Nat
758        Methods* **11**, 1144-1146, doi:10.1038/nmeth.3103 (2014).

759  49    Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.
760        CheckM: assessing the quality of microbial genomes recovered from isolates, single
761        cells, and metagenomes. *Genome Res* **25**, 1043-1055, doi:10.1101/gr.186072.114
762        (2015).

763  50    Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using
764        MinHash. *Genome Biol* **17**, 132, doi:10.1186/s13059-016-0997-x (2016).

765  51    Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and
766        accurate genomic comparisons that enables improved genome recovery from
767        metagenomes through de-replication. *ISME J* **11**, 2864-2868,
768        doi:10.1038/ismej.2017.126 (2017).

769  52    Jain, C., Rodriguez, R. L., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
770        throughput ANI analysis of 90K prokaryotic genomes reveals clear species
771        boundaries. *Nat Commun* **9**, 5114, doi:10.1038/s41467-018-07641-9 (2018).

772  53    Marcais, G. *et al.* MUMmer4: A fast and versatile genome alignment system. *PLoS
773        Comput Biol* **14**, e1005944, doi:10.1371/journal.pcbi.1005944 (2018).

774  54    Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome
775        Biol* **5**, R12, doi:10.1186/gb-2004-5-2-r12 (2004).

776  55    Nguyen, L. T., Schmidt, H. A., von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and
777        effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol
778        Biol Evol* **32**, 268-274, doi:10.1093/molbev/msu300 (2015).

779  56    Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically
780        annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic
781        Acids Research* **47**, D309-D314, doi:10.1093/nar/gky1085 (2018).

782  57    Ashburner, M. *et al.* Gene ontology: tool for the unification of biology. The Gene
783        Ontology Consortium. *Nat Genet* **25**, 25-29, doi:10.1038/75556 (2000).

784    58    Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a
785          reference resource for gene and protein annotation. *Nucleic Acids Res* **44**, D457-462,
786          doi:10.1093/nar/gkv1070 (2016).

787    59    Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL): an online tool for phylogenetic
788          tree    display    and    annotation.    *Bioinformatics*    **23**,    127-128,
789          doi:10.1093/bioinformatics/btl529 (2007).

790    60    Lokmer, A. *et al.* Use of shotgun metagenomics for the identification of protozoa in
791          the gut microbiota of healthy individuals from worldwide populations with various
792          industrialization levels. *PLoS One* **14**, e0211139, doi:10.1371/journal.pone.0211139
793          (2019).

794    61    Schmidt, T. S. *et al.* Extensive transmission of microbes along the gastrointestinal
795          tract. *Elife* **8**, doi:10.7554/eLife.42693 (2019).

796    62    Heintz-Buschart, A. *et al.* Integrated multi-omics of the human gut microbiome in a
797          case   study   of   familial   type   1   diabetes.   *Nat   Microbiol* **2**,   16180,
798          doi:10.1038/nmicrobiol.2016.180 (2016).

799    63    Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
800          DIAMOND. *Nat Methods* **12**, 59-60, doi:10.1038/nmeth.3176 (2015).

801    64    Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
802          alignment search tool. *J Mol Biol* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2
803          (1990).

804

## Competing interests

806    The authors declare no competing interests.

807

## Author contributions

809    CYK and IL conceived this study. CYK and ML constructed the catalog and performed

810    bioinformatics analysis. SY constructed the web server. KK, DY, and HRK organized the

811    study cohorts and provided the fecal samples. IL supervised the project. CYK, ML, and IL

812    wrote the manuscript.

813

## Acknowledgments

818

## 819 Figure legends

820 **Fig. 1 | Effect of sequencing depth on *de novo* genome assembly. a,** Sequencing depth of
821 samples from Korea, Japan, and India. Red data points, nine samples used for the generation
822 of simulated samples for different sequencing depths. **b,** Total read length of samples from
823 Korea, Japan, and India. **c,** The average number of HQ and MQ genomes (left axis) and the
824 proportion of HQ genomes (right axis) from nine samples. **d,e,** Completeness (d) and N50 (e)
825 of assembled genomes from lower sequencing depth (left box of each column) and greater
826 sequencing depth (right box of each column). **f,** The number of the assembled genomes from
827 Korea, Japan, and India. **g,** Total number of the assembled genomes from Korea, Japan, and
828 India, and genome assembly yields. **h,** The relative abundance of 224 Korea-specific, 338
829 Japan-specific, and 18 India-specific assembled genomes in independent fecal samples from
830 the US (n = 926). *P*-values were calculated by two-sided Mann–Whitney U test (**: $P < 0.01$;
831 ***: $P < 0.001$).

832 **Fig. 2 | Phylogenetic tree of 5,386 representative genomes of prokaryotic species from
833 the human gut contained in the HRGM.** Maximum-likelihood phylogenetic tree
834 reconstructed from 120 bacterial marker genes (**Methods**). Representative genomes were
835 annotated by their isolated genome availability (1st layer from the inside), phylum
836 classification (2nd layer), whether they were from UHGG or assembled from KIJ samples
837 (3rd layer), 16S rRNA sequence availability (4th layer), and genome completeness (the
838 outermost layer). Red branches represent 410 genomes from the *Bacteroidaceae* family that
839 are enriched in the representative genome set updated by including KIJ samples.

840 **Fig. 3 | SNV density analysis of the relationship between within-species variation and
841 speciation of gut microbes. a,** The number of SNVs per kb pair of the aligned region. SNV
842 density is summarized for each phylum. Boxes are sorted by the median. Arc, archaeal
843 phylum. **b,** The phylogenetic tree for Actimobacteriota phylum. Inside annotation indicates
844 the *Collinsella* genus, divided into *Collinsella* with modest phylogenetic dispersion (MD

845   *Collinsella*, Red) and *Collinsella* with high phylogenetic dispersion (HD *Collinsella*, Orange).

846   Black annotations in the outer circle represent *Collinsella aerofaciens*, *Collinsella*

847   *aerofaciens_A*, *Collinsella aerofaciens_E*, and *Collinsella aerofaciens_F*, according to the

848   GTDB-TK annotation. **c,** GTDB-TK based taxonomic annotation of MD *Collinsella* and HD

849   *Collinsella*. **d,** SNV density of HD *Collinsella*, MD *Collinsella*, Non-*collinsella*

850   actinobacteriota, and other species. **e,** Scatter plot analysis of SNV density and average

851   phylogenetic distance to the five nearest species of each representative species. Orange points

852   denote species of HD *Collinsella* and black points represent other species. **f,** Comparison of

853   SNV density between the top 10% and bottom 90% species sorted from the lowest average

854   phylogenetic distance to the five nearest species. Statistical significance was calculated by

855   two-sided Mann–Whitney U test (n.s.: not significant; *: $P < 0.05$; ***: $P < 0.001$).

856   **Fig. 4 | Effect of HRGM on taxonomic and functional classification of sequencing reads.**
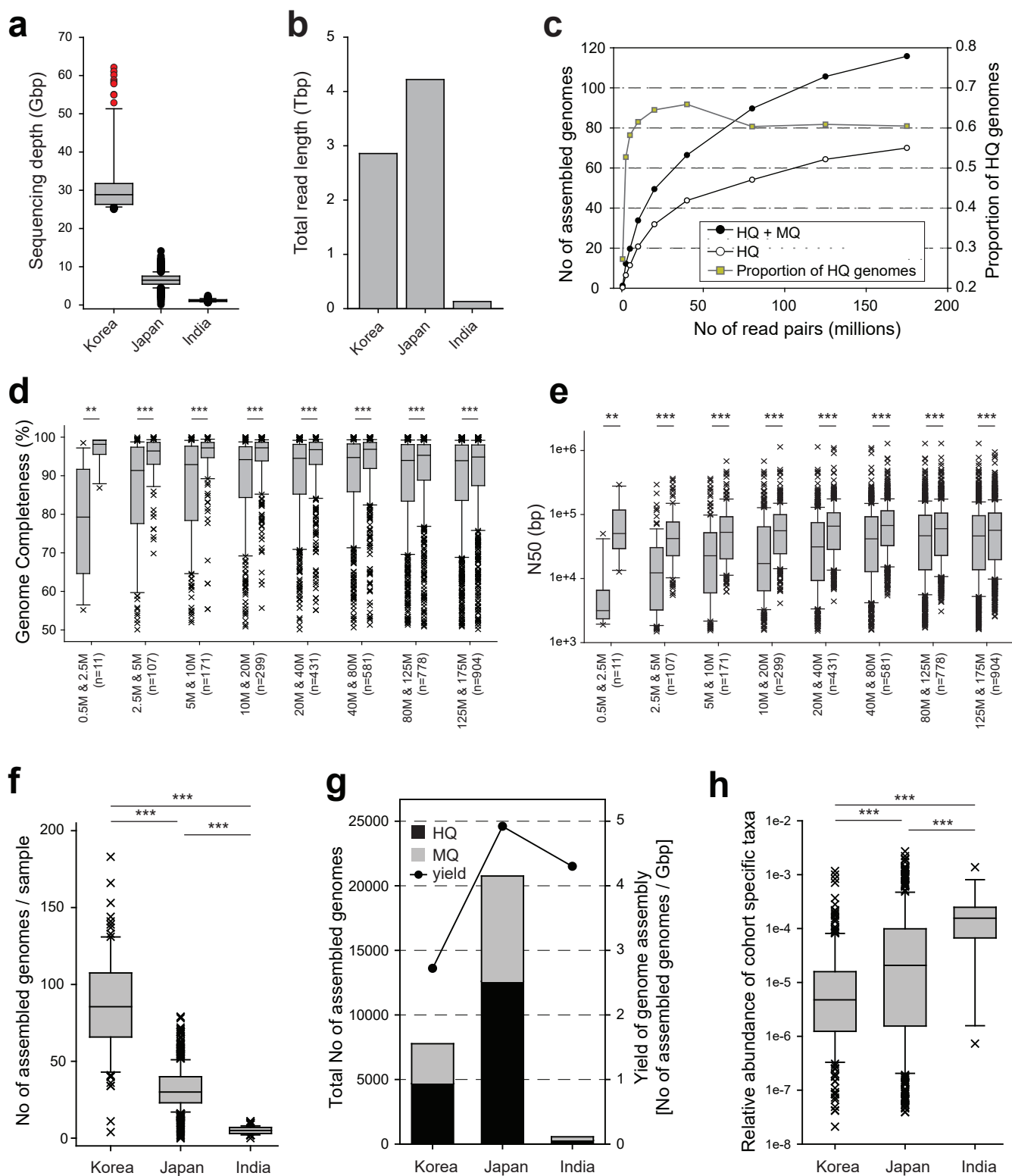
857   **a,** Proportion of taxonomically classified sequencing reads of WMS data from four different

858   populations. The significance of the improvement was calculated by Wilcoxon signed-rank

859   test. Brown–Forsythe test was used to evaluate the decrease of variance. **b,c,** Percent

860   improvement of the read classification proportion in HRGM-based database compared with

861   the standard database (b) and the UHGG-based database (c). **d,** The number of reads aligned

862   to the UHGP-95 and HRGM-95 protein catalogs. Statistical significance was calculated by
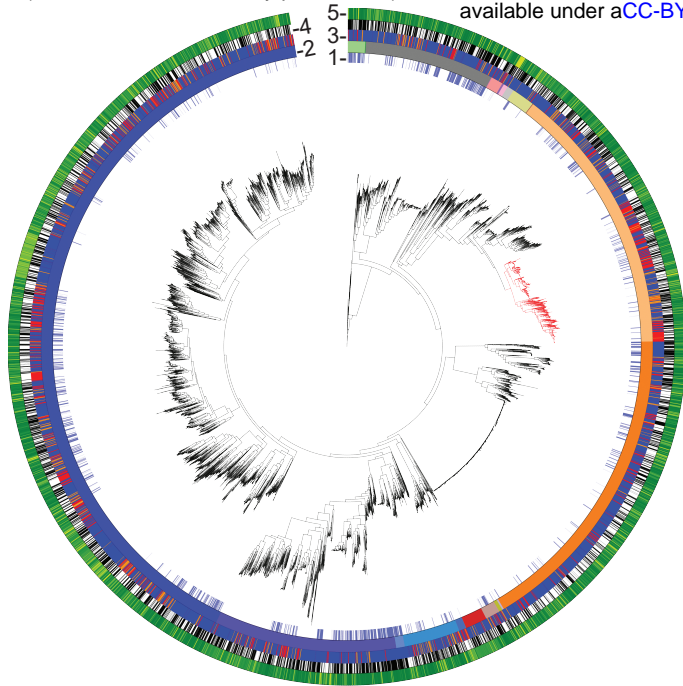
863   using Wilcoxon signed-rank test.

864   **Fig. 5 | Effect of sequencing depth on the reliability of taxonomic profiles. a**, The

865   distribution of taxonomic features over different mean relative abundances. **b**, The

866   cumulative proportion of taxonomic features at different thresholds of mean relative

867   abundance. **c,d,** Pearson correlation coefficient (*PCC*) (c) and Spearman correlation

868   coefficient (*SCC*) (d) of the taxonomic profiles at the given sequencing depth and 80M

869   fragments. The x-axis (the mean relative abundance threshold) indicates the upper boundary

870   of the mean relative abundance.

871   **Fig. 6 | Landscape of cross-reactivity potential of gut prokaryotic genomes. a,** The

872   number of genes and autoimmune epitope sequence-containing genes (ECG) in 5,414

873   genomes of species representatives. Red and orange points, species with the top 1% and 5%

874   ECG per gene, respectively. **b,** Volcano plot of the enrichment of ECG density. Taxonomic

875   clades with positive log2 fold-change and $P < 1e–5$ are highlighted with different colors.

29

876    Taxonomic clades denoted by the same color have an inclusive relationship (e.g.,

877    *g_Prevotella* belongs to *f_Bacteroidaceae*), with the exception of p_Bacteroidota,

878    c_Bacteroidia, and o_Bacteroidales. The first character of each clade name indicates the

879    taxonomic levels (p: phylum; c: class; o: order; f: family; and g: genus). **c,** The red-

880    highlighted area from (b). **d,** Maximum-likelihood phylogenetic tree with taxonomic

881    annotations of clades with high ECG density. The first layer represents clades with the top 1%

882    (red) and 5% (orange) ECG density [annotations and color designations are the same as in

883    (**a**)]. The second and third layers represent enriched taxonomic clades in the volcano plot

884    [taxonomic annotations and color designations are the same as in (b) and (c)]. The second

885    layer represents above-genus level annotations. The third layer represents genus-level
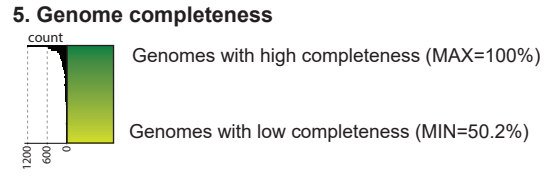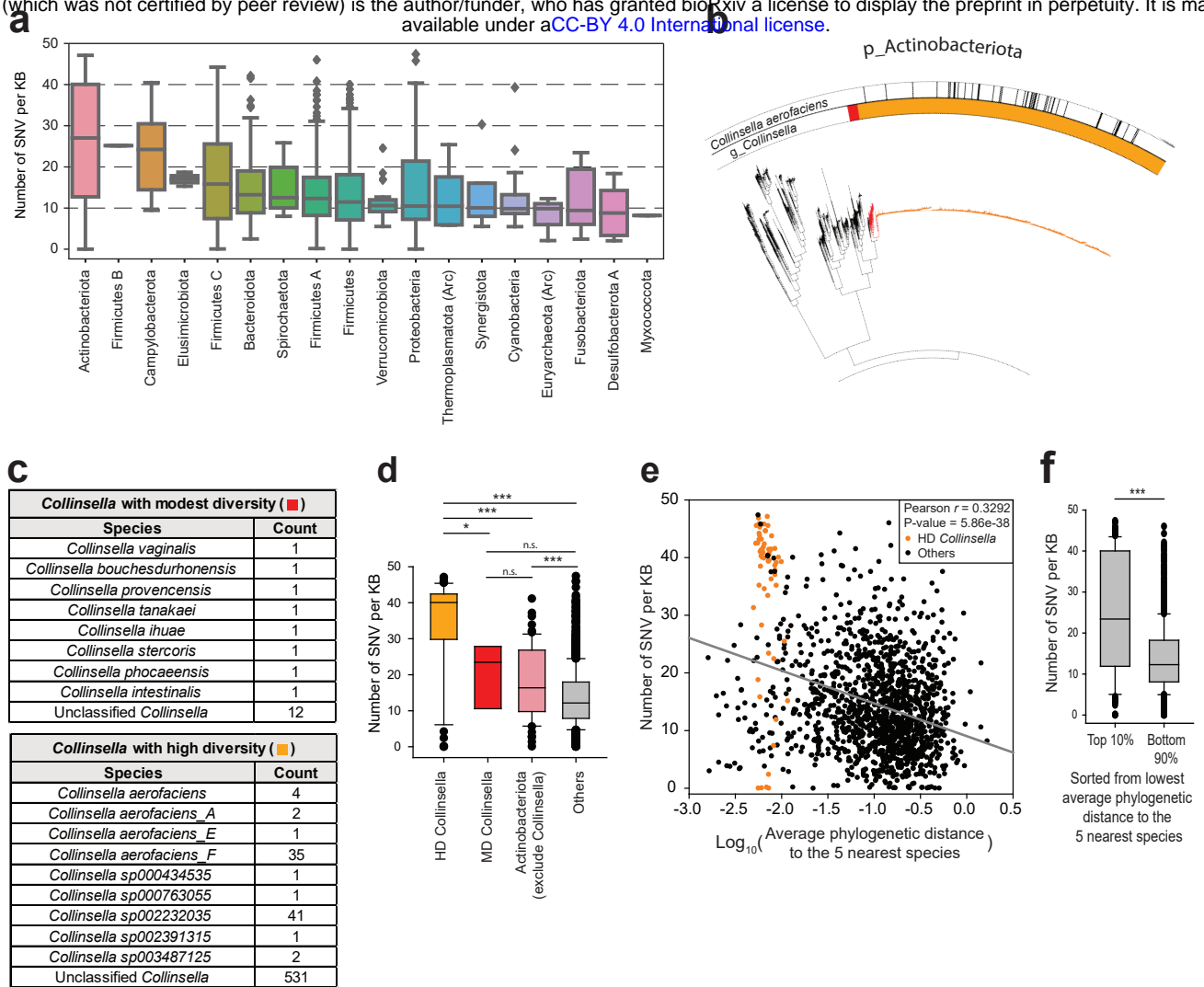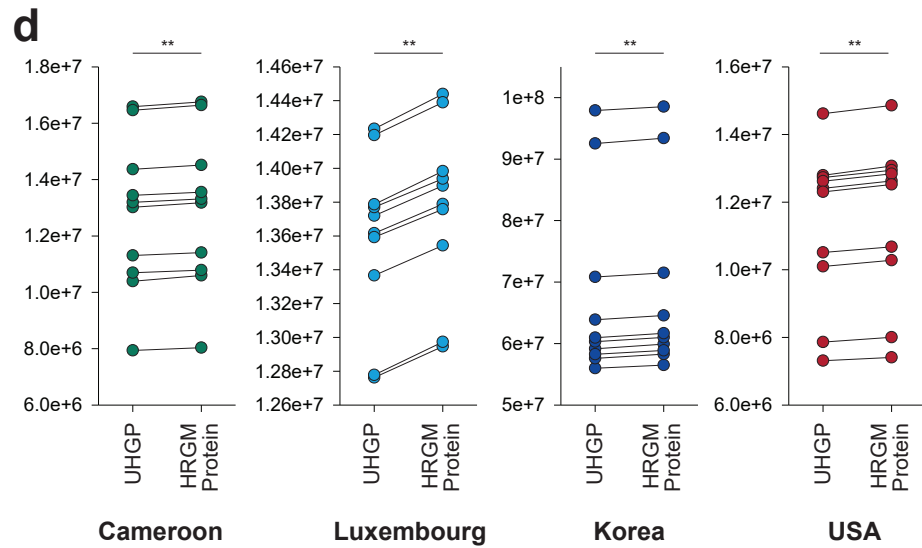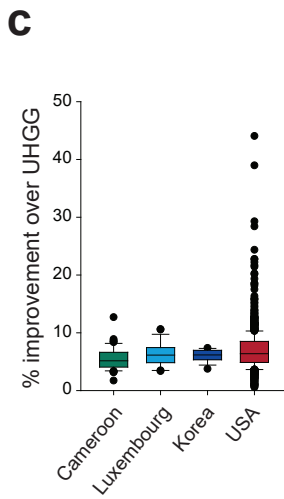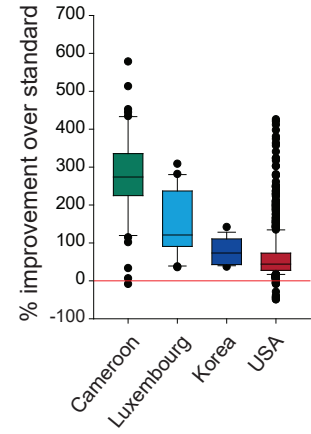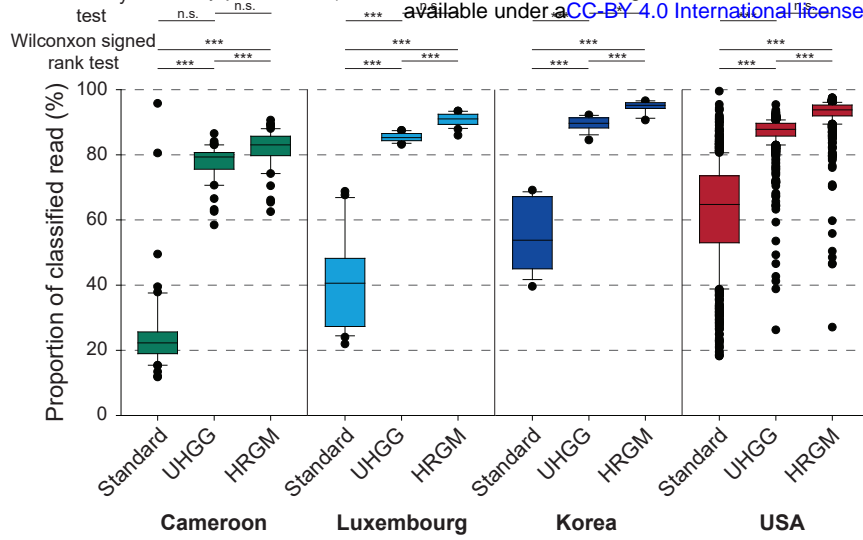
886    taxonomic clades.

887

**1. Isolated genome availability**

- Isolated genomes available (n=893,16.5%)
- No isolated genome available (only MAGs) (n=4,521,83.5%)

**2. Phylum classification**

- Actinobacteriota
- Bacteroidota
- Bdellovibrionota
- Campylobacteriota
- Cyanobacteria
- Desulfobacterota_A
- Elusimicrobiota
- Eremiobacterota
- Fibrobacterota
- Firmicutes
- Firmicutes_A
- Firmicutes_B
- Firmicutes_C
- Firmicutes_G
- Firmicutes_I
- Fusobacteriota
- Myxococcota
- Patescibacteria
- Proteobacteria
- Spirochaetota
- Synergistota
- Verrucomicrobiota

**3. Genomes from UHGG or assembled from KIJ samples?**

- Asssembled from KIJ samples only (n=780,14.4%)
- From UHGG & assembled from KIJ samples (n=580,10.7%)
- From UHGG only (n=4,054, 74.9%)

**4. 16S rRNA sequence availability**

- 16S rRNA sequence available (n=2,542, 47.0%)
- No 16S rRNA sequence available (n=2,872, 53.0%)

**5. Genome completeness**

count

Genomes with high completeness (MAX=100%)

Genomes with low completeness (MIN=50.2%)

1200 600 0

**a**



**b**

p_Actinobacteriota



**c**

**Collinsella with modest diversity ( ■ )**

| Species | Count |
|---|---|
| Collinsella vaginalis | 1 |
| Collinsella bouchesdurhonensis | 1 |
| Collinsella provencensis | 1 |
| Collinsella tanakaei | 1 |
| Collinsella ihuae | 1 |
| Collinsella stercoris | 1 |
| Collinsella phocaeensis | 1 |
| Collinsella intestinalis | 1 |
| Unclassified Collinsella | 12 |

**Collinsella with high diversity ( ■ )**

| Species | Count |
|---|---|
| Collinsella aerofaciens | 4 |
| Collinsella aerofaciens_A | 2 |
| Collinsella aerofaciens_E | 1 |
| Collinsella aerofaciens_F | 35 |
| Collinsella sp000434535 | 1 |
| Collinsella sp000763055 | 1 |
| Collinsella sp002232035 | 41 |
| Collinsella sp002391315 | 1 |
| Collinsella sp003487125 | 2 |
| Unclassified Collinsella | 531 |

**d**



**e**



Pearson $r$ = 0.3292
P-value = 5.86e-38
● HD Collinsella
● Others

$Log_{10}($ Average phylogenetic distance to the 5 nearest species $)$

**f**



Sorted from lowest average phylogenetic distance to the 5 nearest species

**a**



**b**

| Threshold | Feature count | Percentile |
|-----------|---------------|------------|
| < 1e-07 | 24 | 0.34% |
| < 1e-06 | 976 | 13.92% |
| < 1e-05 | 3838 | 54.72% |
| < 1e-04 | 5942 | 84.72% |
| < 1e-03 | 6649 | 94.80% |
| < 1e-02 | 6942 | 98.99% |
| < 1e-01 | 7001 | 99.83% |
| < 1 | All = 7013 | 100% |

**c**



**d**