Open access • Posted Content • DOI:10.1101/528737

# An integrated metagenome catalog reveals novel insights into the murine gut microbiome — Source link ↗

Till Robin Lesker, Abilash Chakravarthy, Eric J. C. Gálvez, Ilias Lagkouvardos ...+7 more authors

**Institutions:** Technische Universität München, Max Planck Society, University of Kiel, RWTH Aachen University ...+1 more institutions

Related papers:

- An Integrated Metagenome Catalog Reveals New Insights into the Murine Gut Microbiome

- An Experimental Metagenome Data Management and Analysis System

- PanFP: pangenome-based functional profiles for microbial communities

- Accurate Annotation of Microbial Metagenomic Genes and Identification of Core Sets

- HumGut: a comprehensive human gut prokaryotic genomes collection filtered by metagenome data

1   **An integrated metagenome catalog reveals novel insights into the murine gut**

2   **microbiome**

3

4   Till Robin Lesker[1], Abilash Chakravarthy[1], Eric. J.C. Gálvez[1], Ilias Lagkouvardos[2],

5   John F. Baines[3,4], Thomas Clavel[2,5], Alexander Sczyrba[6,7], Alice C. McHardy[7,8], Till

6   Strowig[1,9,10]

7

8   Affiliations:

9   1: Department of Microbial Immune Regulation, Helmholtz Centre for Infection

10  Research, Braunschweig, Germany.

11  2: ZIEL Institute for Food and Health, Technical University of Munich, Freising,

12  Germany

13  3: Max Planck Institute for Evolutionary Biology, Plön, Germany.

14  4: Institute for Experimental Medicine, Kiel University, Kiel, Germany.

15  5: Functional Microbiome Research Group, Institute of Medical Microbiology, RWTH

16  University Hospital, Aachen, Germany

17  6: Center for Biotechnology, Bielefeld University, Bielefeld, Germany.

18  7: Department of Computational Biology of Infection Research, Helmholtz Centre for

19  Infection Research, Braunschweig, Germany.

20  8: Braunschweig Integrated Centre of Systems Biology, Braunschweig, Germany

21  9: Hanover Medical School, Hannover, Germany

22  10: RESIST, Cluster of Excellence 2155, Hanover Medical School, Hanover, Germany

23

24  Corresponding author: Till Strowig, till.strowig@helmholtz-hzi.de

25

27  **Abstract**

28  The vast complexity of host-associated microbial ecosystems requires generation of

29  host-specific gene catalogs to survey the functions and diversity of these communities.

30  We generated a comprehensive resource, the integrated mouse gut metagenome

31  catalog (iMGMC), comprising 4.6 million unique genes and 660 high-quality

32  metagenome-assembled genomes (MAGs) linked to reconstructed full-length 16S

33  rRNA gene sequences. iMGMC enables unprecedented coverage and taxonomic

34  resolution, i.e. more than 89% of the identified taxa are not represented in any other

35  databases. The tool (github.com/tillrobin/iMGMC) allowed characterizing the diversity

36  and functions of prevalent and previously unknown microbial community members

37  along the gastrointestinal tract. Moreover, we show that integration of MAGs and 16S

38  rRNA gene data allows a more accurate prediction of functional profiles of communities

39  than based on 16S rRNA amplicons alone. Integrated gene catalogs such as iMGMC

40  are needed to enhance the resolution of numerous existing and future sequencing-

41  based studies.

42

43  **Introduction:**

44  The gut microbiota is a dynamic and highly diverse microbial ecosystem that impacts

45  the hosts physiology[1]. Culture-independent methods such as high-throughput

46  sequencing have revolutionized experimental approaches to characterize and

47  investigate these communities. Gene catalogs facilitate taxonomic and functional

48  annotation of sequencing data, thereby maximizing insights gained from short-reads[2–

49  5]. Moreover, they can provide higher resolution than less specific resources such as

50  GenBank by including valuable metadata such as environment-specific variables.

51  Typically, generation of reference gene catalogs involves sample-specific assembly,

52  prediction of genes and dataset-wide clustering of gene entries to reduce redundancy.

53  However, this approach results in reduced taxonomic resolution of gene entries, first

54  due to clustering of highly related but distinct genes and second due to the lack of high-

55  resolution taxonomic information for gene entries, which can be best obtained from

56  marker genes, such as 16S rRNA genes for which large reference collections exist.

57  Here we present a novel approach and corresponding computational workflow to

58  construct integrated gene catalogs, resulting in a significant improvement of the

59  taxonomic resolution of gene entries and providing valuable additional information

60  such as linking genes to metagenome-assembled genomes (MAGs) and

61  reconstructed full-length 16S rRNA genes. We applied this approach to construct an

62  integrated mouse gut metagenome catalog (iMGMC) combining existing and newly

63  sequenced metagenomic data. We chose this ecosystem as the mouse serves as

64  foremost experimental model system to study microbiota-modulated human diseases,

65  but the use of currently existing human gut gene catalogs is precluded due to the

66  substantial differences in bacterial species and genes present in mice[6].

67

68  **Results:**

69  ***Construction of the integrated mouse gut metagenome catalog (iMGMC)***

70  Pioneering work by others resulted in the construction of several gene catalogs,

71  including a microbiome gene catalog from the mouse gut (hereon referred to as

72  MGCv1 ) comprising 2.6 million non-redundant genes[4]. We developed a bioinformatic

73    workflow that combines a global assembly strategy with binning of contigs to putative

74    MAGs and with innovative linking of reconstructed 16S rRNA gene sequences to these

75    MAGs (Figure 1A). This "All-in-One" assembly approach together with the subsequent

76    binning enables maintaining complex information such as distribution of distinct contigs

77    and bins over a large number of samples. We applied this approach to a previously

78    published set of sequencing data included in MGCv1 (n = 190 mouse fecal samples)

79    and increased the biological diversity by incorporating novel metagenomic data for 108

80    additional intestinal samples from a large number of commercial mouse providers and

81    wild mice, including different gastrointestinal locations (see Table S1). This selection

82    was based on the previous notion that the source of experimental mice and anatomic

83    niches contribute to the variability between murine microbiome to a higher extent than

84    other factors such as diet, genotype, housing laboratories or gender[4]. As a first step in

85    the construction of iMG2C, 1.3 Tbp from 298 metagenomic sequencing libraries were

86    assembled using Megahit[7] in an "All-in-One" approach, resulting in 1.2 million contigs

87    of length greater than 1000bp, with a total assembly size of 4.5 Gbp. Next, genes were

88    identified with MetaGeneMark[8], resulting in 4.6 million open reading frames (ORFs) of

89    length greater than 100 bp, compared to 2.6 million ORFs in the MGCv1 (+77%)

90    (Figure 1B). We tested the redundancy of these ORFs by clustering them with CD-Hit

91    (95% identity at 90% coverage)[9], which resulted in a reduction of only 2% of ORFs (n

92    = 99,670) (data not shown). We considered this negligible compared to the 89%

93    reduction in MGCv1 [4]. Subsequently, contigs were binned using MetaBat[10], resulting

94    in 1,462 bins greater than 200 kbp (containing 87% of iMGMC entries). Subsequently,

95    we defined 660 bins encoding 40% of all iMGMC entries as MAGs, based on the

96    presence of established sets of bacterial marker genes and a quality threshold ≥80%

97    (Figure 1C) [11]. Notably, MGCv1 did not provide MAGs, as sample-specific assemblies

98    were used, but rather less specific information referred to as "co-abundance groups"

99    (CAGs), containing at least 700 genes. Comparison of the numbers of CAGs and

100   genes in CAGs between iMGMC and MGCv1 revealed large increases in our resource

101   (1,217 vs. 541 CAGs, 81% vs. 40% of genes, respectively) (Figure 1B).

102   In addition to reconstructing bins including MAGs, we also assembled 16S rRNA

103   genes, using the following approach that overcomes the limitation that 16S rRNA

104   genes are typically not efficiently recovered in standard assemblies, due to their highly

105   conserved regions[12]: Using RAMBL[13], we reconstructed 1,323 full-length, unique 16S

106   rRNA gene sequences, a number similar to the number of genomes (n=1,068)

107   predicted based on the presence of 139 distinct marker genes in the iMGMC assembly

108   using Anvi'o (Figure 1E)[14]. We postulated that linking 16S rRNA genes to bins and

109   MAGs after assembly would allow efficient integration of these complementary pieces

110    of information, thereby improving the taxonomic assignment of MAGs. However, no

111    high-throughput method currently exists for creating such links. Hence, we developed

112    an integrated score combining mapping- and correlation-based associations to assign

113    a 16S rRNA gene sequence to each bin and MAG (Figure 1F and S1). Briefly, we first

114    identified all contigs containing reconstructed 16S rRNA gene sequences via BlastN

115    [15]. Then, we searched for paired-end reads in which one read mapped to a

116    reconstructed 16S rRNA gene sequence and the other to a contig. Finally, we

117    remapped all libraries to the 1,462 bins and the 1,323 16S rRNA gene sequences to

118    determine their relative abundances across all samples and used this data to estimate

119    correlations between bins and 16S rRNA gene sequences using an abundance co-

120    variance strategy[16]. This individual information was finally integrated using a novel

121    approach (see Methods for details) to assign the reconstructed 16S rRNA genes to

122    bins.

123

124    ***Evaluation of iMGMC generation***

125    The different steps underlying the construction of iMGMC were evaluated for their

126    efficiency using those MAGs that had a highly related reference genome. These were

127    specifically identified by mapping synthetic reads generated with BBMap from all 9,748

128    bacterial genomes available in the NCBI Assembly database (Version January 2017)

129    against all bins and also the contigs that we were not able to bin (unbinned contigs) in

130    iMGMC (see Methods for details). After read mapping, we evaluated the distribution of

131    these genomes in our assembly and identified 57 genomes, which were recovered at

132    least by 50% within binned and unbinned contigs. For these genomes, we recovered

133    on average 79 ± 11% (mean ± s.d.) in our assembly, from which 78 ± 19% were found

134    in the respective best/largest bin, while only 13 ± 17% were found in unbinned contigs

135    (Figure 1G and S2). Thus, we considered our "All-in-one" assembly as good as other

136    assembly strategies employed for large-scale MAG reconstruction[17]. The number of

137    MAGs (n=660) would even be higher when using a quality threshold from an already

138    published study (n=818, quality(CheckM): Completeness – 5x contamination $\geq$ 50%)[17].

139    We also evaluated the utility of the "All-in-one" assembly approach for another large

140    dataset by processing metagenomic sequencing data from the pig microbiome. From

141    287 fecal samples (1,758 Gb) used to construct a previous reference gene catalog[5],

142    we obtained 12.2 Mio ORFs and 1,050 MAGs, representing a 58%- and 45 %-

143    increase, respectively, compared with the original work (data not shown).

144    The MAG/16S rRNA gene pairs were evaluated using MAGs with linked 16S rRNA

145    gene sequences for which reference genomes exist. Specifically, we identified

146    genomes found in our assembly and the respective bins, followed by comparison of

147  the known 16S rRNA gene sequences to the correspondingly predicted 16S rRNA

148  gene sequences (Figure S3) (see Methods for detail). From the 47 identified genomes

149  and respective bins, 28 agreed perfectly (100% sequence identity) between known

150  and linked 16S rRNA gene, with an additional 7 matching taxonomic assignment down

151  to the genus level. The remaining 12 genomes and bins disagreed at varying

152  taxonomic levels (Figure 1H and S3). Statistical assessment of these results supported

153  that our approach i) did not require 16S rRNA gene sequences within a MAG to

154  successfully perform a matching linking and ii) performed better than a random

155  assignment (P=0.074, Pearson's Chi-squared test with Yates' continuity correction).

156  Hence, the proposed novel scoring scheme is with high confidence able to link MAGs

157  and bins to corresponding reconstructed 16S rRNA genes, improving taxonomic

158  resolution, though not in an error-free manner.

159

160  Thus, we created a novel type of resource which i) includes a gene catalogs that

161  outperform previous versions and ii) includes novel information, i.e. MAGs, and 16S

162  rRNA gene sequences, which are linked with each other.

163

164  **iMGMC reveals high prevalence of novel taxa in the mouse gut microbiota**

165  Both metagenomic and cultivation-based studies showed that the gut microbiome of

166  mice compared to human is composed of distinct bacterial species, of which many are

167  yet uncultured and lack genomic information[4,6]. Analysis of our 660 reconstructed

168  MAGs corroborates this notion, revealing that only 72 of them have closely related

169  NCBI assemblies including other MAGs available (ANI > 95%) (Data in Table S1)[18]. A

170  similar observation (only 137 known of 1,050 MAGs in total) was made for MAGs

171  derived from the pig microbiome.

172  To construct a comprehensive phylogenetic tree of the mouse gut microbiota, we

173  assigned MAGs (n=660) and closely related, previously sequenced genomes (n=64)

174  into clusters (Figure 2). In line with previous reports[6,19], our data analysis corroborates

175  that the murine gut microbiome is overall dominated by the two main phyla *Firmicutes*

176  (77% of MAGs / 73% of 16S rRNA gene sequences) and *Bacteroidetes* (14% /

177  18%)(Figure 2 and S4). Notably, *Bacteroidetes* included the second largest MAG

178  cluster, namely the *Bacteroidales* S24-7 group (64% / 49%), recognized as being very

179  abundant in the mouse gut, but for which only three reference genomes are available

180  [6](new Microbiome paper). Strikingly, ≥13 % of MAGs were from phylogenetic groups

181  (up to level of family) that completely lacked reference genomes in public databases

182  (NCBI genomes RefSeq, not other MAGs), such as MAGs assigned to the

183  *Clostridiales*-vadinBB660 group (n= 70) and *Mollicutes* RF9 (n=14) (Figure 2).

5

184  Unsupervised clustering of MAG according to their functional potential (Figure S5)
185  demonstrated that distinct taxonomic clusters such as *Clostridiales*-vadinBB660 group
186  or the *Bacteroidales* S24-7 group represent functionally distinct microbes within the
187  mouse microbiome (Figure S7) (new Microbiome paper).
188  Many additional undescribed bacteria were also identified after comparing the
189  reconstructed 16S rRNA gene sequences to members of "16S ribosomal RNA
190  (Bacteria and Archaea)" at the NCBI-database, with only 164 of 1,323 (12%) having at
191  least a 97% identical match. A large fraction of these sequences were neither found in
192  the SILVA SSU Ref v. 128 database (99% ident: 72% new, 97% ident: 45% new) nor
193  in a recent 16S rRNA database established by target-specific environment
194  sequencing[20] (99% ident: 98% new, 97% ident: 93% new). Notably, while the MAGs
195  represent a large fraction of the phylogenetic tree of the bacteria present in the mouse
196  gut, several taxonomic groups were represented by 16S rRNA gene sequences, but
197  underrepresented by MAGs, such as the family of *Prevotellaceae* (49 16S rRNA gene
198  sequences / 3 MAGs), the class of *Bacilli* (81/10) as well as the phyla of *Proteobacteria*
199  (67/24) and *Actinobacteria* (78/22) (Figure S4). Thus, our analysis identified taxonomic
200  groups that are interesting novel targets for cultivation-dependent and -independent
201  studies to extend our understanding of microbiome-modulated phenotypes in mouse
202  models.
203
204  **Improved functional prediction via MAG/16S rRNA gene links in iMGMC**
205  The establishment of databases of microbial reference genomes has spurred the
206  development of computational approaches to simulate the functional profiles of
207  metagenomes based on marker gene datasets such as 16S rRNA amplicon
208  profiles[21,22]. However, the power of these approaches depends on the availability of
209  sequenced microbial genomes from the respective environments to perform
210  satisfactorily. Because of the existence of numerous bacterial species within murine
211  gut communities that lack reference genomes, we hypothesized that the default
212  PICRUSt-based predictions of mouse-associated metagenome functions are limited[21].
213  Thus, we constructed a mouse-optimized PICRUST version, employing the original
214  PICRUSt algorithm in conjunction with the iMGMC data. Specifically, we used the
215  MAGs with unique linked 16S rRNA sequences (n=484), as well as the 1,322 16S
216  rRNA sequences from the iMGMC to create an extended genome resource for
217  PICRUSt (PICRUSt-iMGMC) (Figure 3A, see methods for details). Comparison of
218  Kegg Ortholog (KO) functional profiles predicted by the default and extended PICRUSt
219  approach using 16S rRNA amplicon data from different gastrointestinal sites (n=50)
220  for the corresponding shotgun metagenomic libraries (WGS) demonstrated a higher

221  correlation to the WGS-based KO profiles for PICRUSt-iMGMC than PICRUSt-default

222  predicted profiles (Pearson: 0.84 vs 0.68, +23%, Spearman: 0.84 vs 0.70, 21%)

223  (Figure 3B and C). The highest correlations were observed for samples from the colon

224  (Pearson: 0.86 vs. 0.67, Spearman: 0.87 vs. 0.72) (Figure S6). Similar improvements

225  were obtained with distinct datasets not used for the construction of the catalog (Figure

226  S7). The improved correlation of PICRUST-iMGMC largely derived from increased

227  sensitivity, i.e. "true positive rates", rather than decreased "false positive rates",

228  enabling the prediction of functionalities otherwise lost (Figure 3D and E). Even when

229  mapping WGS data to the KEGG database with DIAMOND[23] instead of to the iMGMC

230  for generation of the KO reference profile, PICRUSt-iMGMC performed better than

231  PICRUSt-default in predicting functional profiles (Figure S6 and S7).–Finally, we

232  evaluated whether combining the information of iMGMC with the genomes available in

233  the KEGG database improved prediction. Strikingly, PICRUSt-iMGMC/KEGG did not

234  perform better and the correlation with WGS data even decreased, suggesting that

235  inclusion of related but divergent genomes reduces prediction accuracy (Figure S6 and

236  S7). Hence, our resource enabled the development of ecosystem-specific PICRUSt

237  models, i.e. optimized for the murine intestinal microbiome, with substantial

238  improvement in the prediction of metagenomic functional profiles.

239

240  **Multi-scale taxonomic assignment of gene entries based on metagenomic**

241  **reconstruction enhances taxonomic resolution in iMGMC**

242  Gene catalogs have foremost been employed to generate functional profiles from short

243  read metagenome surveys of communities. To assess the performance of iMGMC in

244  this respect, we performed read-mapping of sequencing data from three external

245  studies, which were not included in the construction of neither iMGMC nor MGCv1, to

246  both catalogs[24–26]. This revealed an increased number of reads (up to 36%) mapping

247  to the iMGMC, supporting the utility of this new catalog (Figure S8).

248  The taxonomic assignment of entries in classical gene catalogs, specifically after

249  sample-specific assembly and clustering of ORFs by similarity, i.e. 95% identity at 90%

250  coverage in the MGCv1[4], is limited by the ability of algorithms to predict the taxonomic

251  placement based on relatively short ORFs, which has a limited robustness[27]. Taking

252  advantage of the clustering free approach, we annotated each iMGMC entry using the

253  taxonomic information obtained from the respective gene and contig as well as from

254  the bin [28] and the connected MAG/16S rRNA gene sequences, whenever available

255  (Figure 1D). As a result of using longer contigs rather than short ORFs sequences, the

256  relative taxonomic assignment rate improved between 28 and 1,021% at different

257  taxonomic levels (Figure 1D). Notably, many entries were still not assigned to high

258 taxonomic ranks with high confidence, since these approaches are reference-based,
259 and are hampered by the presence of novel and unclassified taxa. Using the MAGs of
260 the iMGMC resource, we could assign up to 40% of mapped reads of three external
261 datasets to MAGs (Figure S8), facilitating the identification of specific bacterial taxa,
262 allowing improved functional analysis by providing information of the genomic context
263 of genes, or of bacterial interaction networks identified by covarying abundances
264 across samples. For instance, the analysis of previously generated shotgun
265 metagenomic data from mice subjected to different experimental diets allowed the
266 retrospective identification of MAG networks rather than gene clusters that show
267 conserved changes in their relative abundance induced by these diets (Figure S9).
268 Hence, future users will be able to utilize in parallel taxonomic information for each
269 gene catalog entry, ranging from well-established methods with lower resolution to
270 innovative methods with enhanced resolution.
271

272 **Provider-specific diversification of the mouse microbiota**
273 Recent studies have demonstrated that the composition of murine microbiomes varies
274 between different providers, mostly via 16S rRNA amplicon sequence analysis [29].
275 However, to which degree laboratory mice share a conserved set of microbes is not
276 known. The presence of a core set of bacteria, based on the detection of 26 CAGs in
277 >95% of mice, was proposed previously [4]. We analyzed the relative abundance of each
278 individual MAG in all samples by remapping all reads from each library to the MAGs,
279 followed by conversion of mapped read counts into relative abundances (see methods
280 for details). Strikingly, this analysis revealed that each mouse line featured a unique
281 combination of MAGs; even mice from different barriers of the same commercial
282 vendor differed (Figure 2). This resulted in substantial differences in the functional
283 potential of the microbiome within each mouse line (Figure S5D, Table S3). Hence, we
284 next quantitatively assessed the distribution of MAGs by determining their prevalence
285 and relative abundance within each provider. Around 10% of MAGs (70/660) were
286 shared by at least half of the providers (> 0.1% relative abundance in at least one
287 individual sample per provider) (Figure 4A). The most prevalent MAG, matching to
288 *Lactobacillus murinus ASF361* (ANI =97%)*, was detected in almost all providers
289 (20/21). Notably, three additional members of the Altered Schaedler Flora (ASF)
290 community, which has been studied as mouse gut model community in the past, as
291 well as only four other previously sequenced bacteria were found in at least 50% of
292 providers, while the remaining 62 (=88%) represent uncultured bacteria. We next
293 analyzed the MAGs shared by at least two thirds of the providers (n=21 MAGs) from
294 which most belonged taxonomically to the *Firmicutes* (n=18), two belonged to the

295     *Bacteroidales* S24-7 group (phylum *Bacteroidetes*, proposed family *Muribaculaceae*)

296     and one was identical to *Mucispirillum schaedleri* (phylum *Deferribacteres*) (Figure

297     4B). Strikingly, the relative abundance of these MAGs revealed large differences

298     between providers (up to 100-fold) suggesting that their respective abundance within

299     each community is strongly influenced by environmental factors. Taking advantage of

300     the link between MAG and 16S rRNA gene sequences, we assessed the global

301     prevalence and relative abundance of the corresponding 16S rRNA gene sequences

302     across all 16S rRNA amplicon datasets deposited in the SRA using the recently

303     established IMNGS database (Figure 4C)[30]. This search revealed that the most

304     prevalent MAG in our study, *Lactobacillus murinus,* is present in 36% of all samples

305     derived from the mouse gut (n=9,496), while being largely absent from the human gut

306     and only detectable in 1.4% of rat gut microbiota samples (1.4% positive) (Table S4).

307     To assess whether the newly reconstructed 16S rRNA gene sequences represented

308     taxa commonly found in mice, we employed IMNGS and queried all 1,323 16S rRNA

309     gene sequences to assess their relative abundance in SRA samples derived from

310     diverse ecosystems (Figure 4D and E). A prevalence of 1% (threshold relative

311     abundance: 0.1%) within at least one of the ecosystems was determined for 739 rRNA

312     gene sequences from which 569 were enriched in the mouse gut, mouse skin, rat gut

313     or human gut. Of these 44% were most prevalent in the mouse gut, with an additional

314     6% being shared with the mouse skin. Other sequences were shared with the rat

315     microbiome (12%) and the human gut microbiome (7%) (Figure 4E). In summary, our

316     large-scale analysis revealed the presence of specific bacteria commonly found in

317     mouse lines but no other gut microbiomes, yet, also a high species-level variability

318     within the murine gut microbiome, which impacts the functional repertoire of the

319     microbiome and potentially thereby the outcome of *in vivo* experiments.

320

321     **Discussion:**

322     Short read-based sequencing studies of microbial ecosystems require suitable

323     reference databases for maximal resolution of taxonomic and functional assignments.

324     Gene catalogs and 16S rRNA gene databases commonly represent separate

325     references for shotgun metagenome and 16S rRNA amplicon sequencing analyses,

326     respectively. To overcome the separation between these types of databases, a novel

327     framework that can serve as i) a valuable resource for the most utilized experimental

328     model for microbiome research, the mouse gut microbiota, and ii) a blue print to

329     generate integrated gene catalogs for less characterized microbial ecosystems was

330     developed.

331     For the establishment of the integrated gene catalog, methods identified to yield

332   optimal results by the CAMI challenge[27], e.g. for assembly of MAGs or binning when

333   dealing with large datasets, were utilized and complemented  with a novel approach

334   linking MAGs and 16S rRNA sequences. The "All-in-One" assembly resulted for the

335   mouse gut microbiome in the reconstruction of a large number of high-quality MAGs,

336   including low abundant community members, representing bacteria that were neither

337   cultured or identified in other high-throughput sequencing studies[17]. Strikingly, for both

338   the mouse and pig gut microbiome, more than 87% of MAGs fell into this category.

339   The *Clostridiales*-vadinBB660 or *Mollicutes* RF9 groups, which were so far only known

340   from 16S rRNA gene sequencing, are examples of functionally distinct and

341   underexplored bacteria frequently occurring in mouse gut microbiomes. Preliminary

342   analysis of assemblies of large datasets from the human gut microbiome suggest that

343   the developed approach also identifies hundreds of novel MAGs (approximately 30%

344   of assembled MAGs), demonstrating the power of this approach even for better

345   characterized ecosystems.

346   Another utility of the integrated gene catalog is the availability of linked MAG-16S rRNA

347   gene pairs, which enables the incorporation of data from large 16S rRNA gene

348   databases such as the IMNGS database encompassing 168,573 short-read datasets

349   (build 1711) thereby allowing large-scale screening for identified MAGs, such as the

350   evaluation of a core microbiome in the mouse gut. The MAG-16S rRNA gene pairs

351   also enabled the development of an ecosystem-optimized version of PICRUSt, which

352   produced gene profiles more closely resembling WGS data. We anticipate this to be

353   widely adapted to predict metagenome profiles based on 16S rRNA amplicon

354   sequencing data and suggest that ecosystem-optimized versions of PICRUSt will be

355   valuable resources.

356   Altogether, the clustering-free construction of gene catalogs together with the

357   reconstruction of a large number of almost complete MAGs through an improved

358   assembly strategy as well as linking to 16S rRNA gene sequences provide a highly-

359   integrated resource for sequencing-based work and will enable future studies to

360   explore the taxonomy, functionality and community structure of the mouse gut and

361   other ecosystems in more depth.

362

11

**Figure 1: Generation and evaluation of the integrated mouse gut metagenome catalog (iMGMC)**

(A) Flowchart displaying the steps and bioinformatics tools (names in brackets) utilized for the generation of the iMGMC. This resource includes genes, metagenome assembled genomes (MAGs), 16S rRNA gene sequences and MAG-16S rRNA gene links.

(B) Comparison of relative and total numbers of gene entries and their association to bins of different completeness between a previous mouse gut gene catalog (MGCv1)[4] and iMGMC. Bins were defined as: i) co-abundance genomes (CAG) if they were larger than >= 200kbp lengths and contained ≥700 ORFs or: ii) MAGs if their quality (marker gene completeness – contamination) as determined by CheckM was ≥ 80%.

(C) Quality determination of individual binned contigs by CheckM by analyzing marker gene completeness and contamination. Box plots display marker gene completeness and contamination of 660 MAGs and 802 CAGs, respectively.

(D) Absolute numbers of gene entries colored according to the lowest possible taxonomic annotation of the ORF, contig or bin. Different taxonomic profilers were employed for classification: ORF: DIAMOND-BlastP; contigs: CAT (Contig annotation tool); bins: GTDBTk

(E) Number of genomes in dataset estimated using a marker gene set containing 139 genes. Each dot represents the copy number of the respective marker gene.

(F) Overview of the methodology to link MAGs to 16S rRNA gene sequences by combining mapping-based and statistical approaches. Resulting linked pairs of MAGs and reconstructed 16S rRNA gene sequences were used together with KEGG annotations for construction of mouse gut specific PICRUSt predictions.

(G) Evaluation of binning by calculating the fraction of recovered RefSeq genomes (threshold ≥ 50 % of genome present in contigs, n=57) in bins.
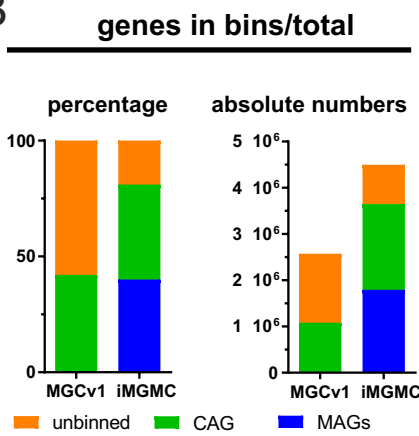
(H) Evaluation of MAG / 16S rRNA gene linking by determining the taxonomic match between predicted and reference 16S rRNA gene sequence for those recovered RefSeq genomes with a MAG / 16S rRNA gene pair (n=47).
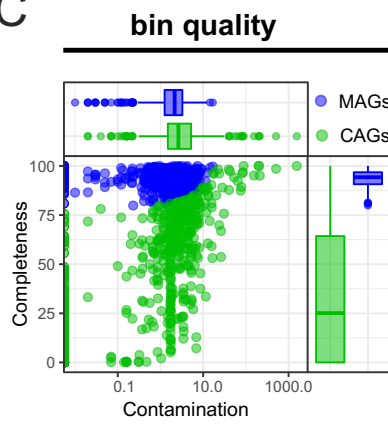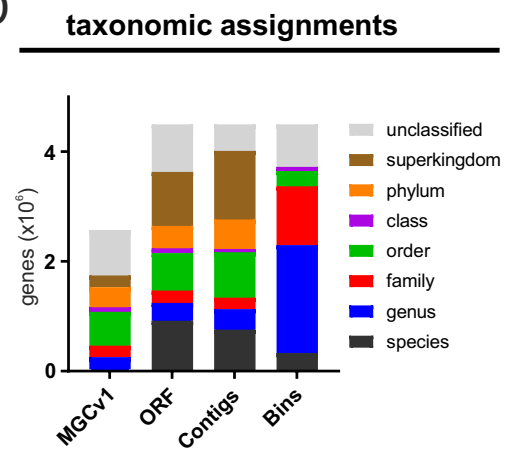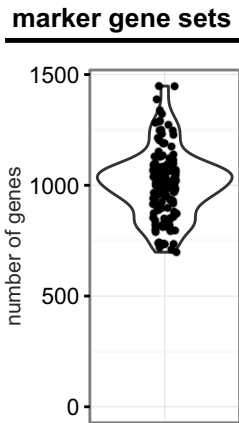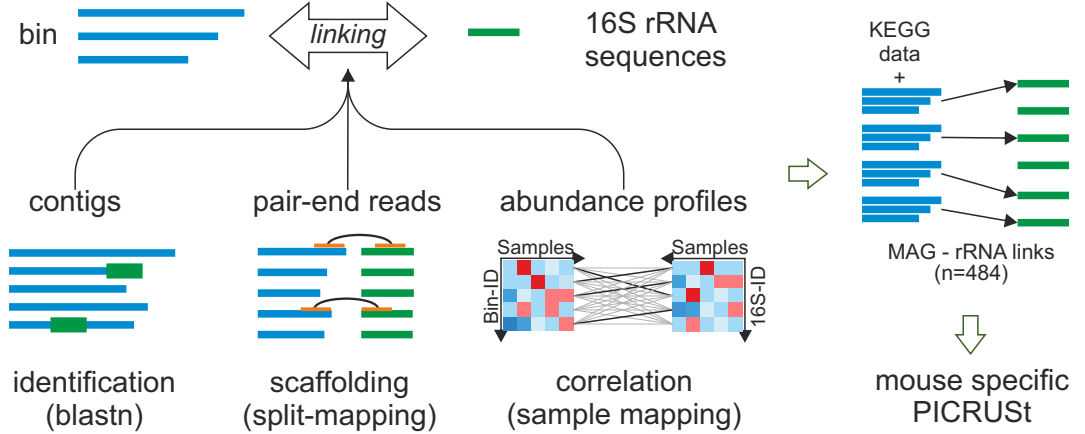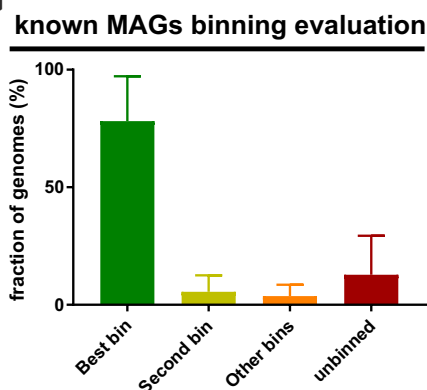
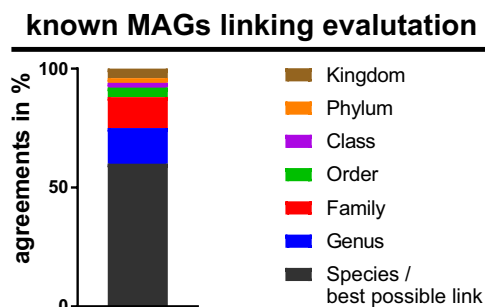See Figures S1, S2 and S3 for more details.

Figure 1

402 **Figure 2: Phylogenetic tree of the 660 MAGs included in the iMGMC**

403 MAGs are shown as triangles and 64 closely related, previously sequenced bacteria

404 used for comparison as stars (genomes from NCBI refSeq with mapping rate >50%

405 coverage). The color of triangles indicates their taxonomic association to different

406 phyla and the size of triangles indicates the mean relative abundance in all iMGMC

407 samples. The tree includes manually curated taxonomic assignments for most MAGs

408 and the names of the taxonomic clusters are displayed in full or abbreviated in the tree.

409 The inner rings show the relative abundance of the 660 MAGs in the 21 investigated

410 mouse providers (threshold: 0.1%). The last three rings visualize the relative

411 abundance of 469 of 660 MAGs at different anatomical sites (threshold: 0.1%, SI: small

412 intestine). The outer bar plots show their respective maximal relative abundance.

413

# MAGs-Tree

Taxa designations
A: *Dorea*
B: *Blautia*
C: *Coprococcus*_1
D: *Lachnoclostridium*
E: *Tyzzerella*_3
F: *Ruminococcaceae*_UCG-013
G: *Ruminococcaceae*_UCG-010
H: *Ruminiclostridium*
I: *Ruminococcus*
J: *Anaerotruncus*
K: *Ruminococcaceae*_UCG-014
L: *Erysipelotrichaceae*
M: *Coriobacteriaceae*
N: *Desulfovibrionaceae*
O: *Alphaproteobacteria*
P: *Gammaproteobacteria*
Q: *Akkermansia*
R: *Odoribacteraceae*
S: *Bacteroides*
T: *Parabacteroides*
U: *Prevotella*

related NCBI-Bacteria (RefSeq)
MGAs (this study)
Phylum-*Bacteroidetes*
Phylum-*Firmicutes/Tenericutes*
Phylum-other
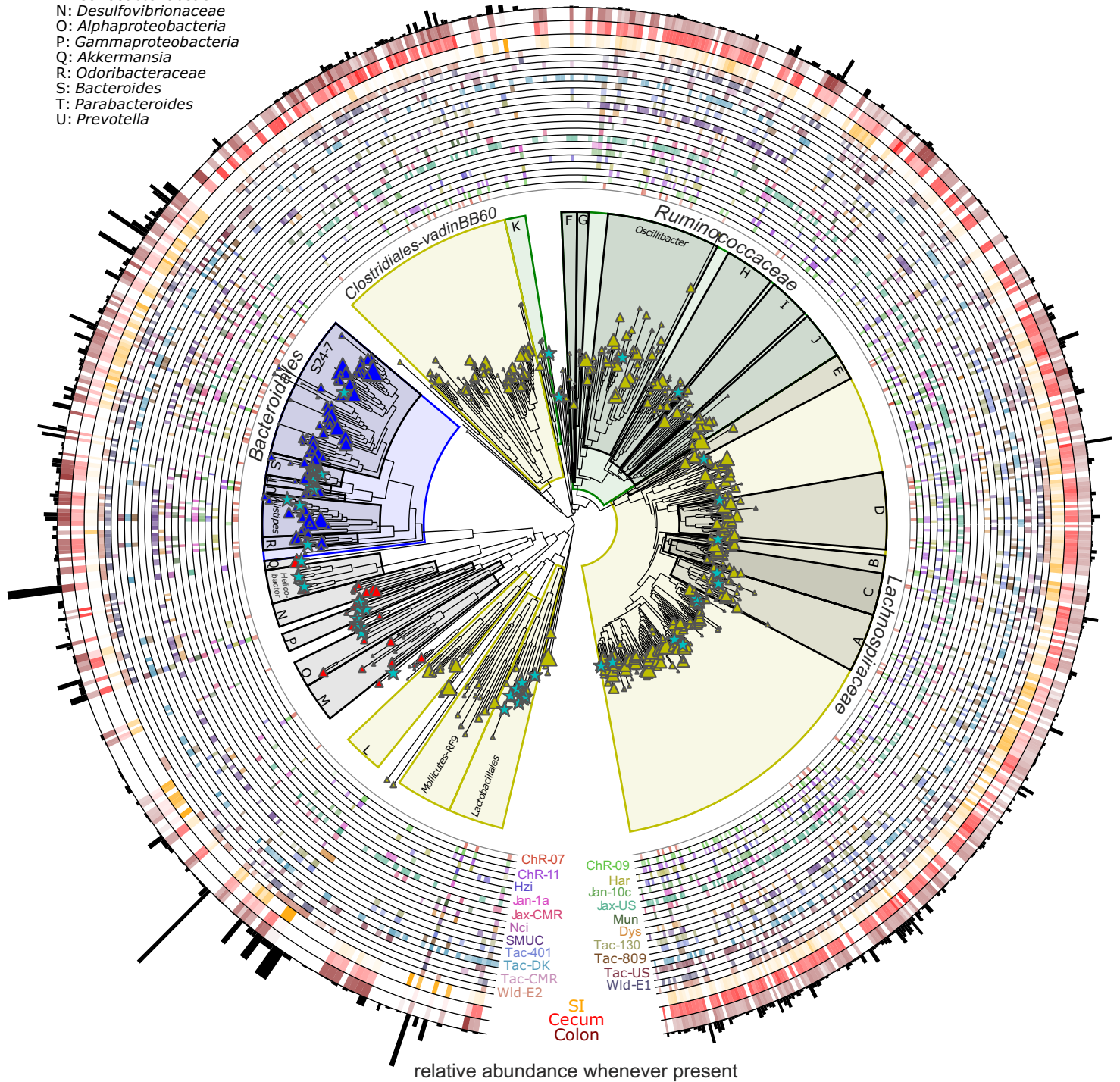


relative abundance whenever present

Figure 2

414    **Figure 3:**

415    **Mouse gut microbiota optimized PICRUSt-iMGMC model**

416    (A) The different PICRUSt workflows used in this study: (I) Default workflow for end-

417    user starting from close reference picked OTUs against the GreenGenes database

418    relying on functional metagenome prediction using precalculated genome predictions

419    files (II) Novel PICRUSt workflow starting from *denovo* picked OTUs and using MAGs

420    with 16S rRNA gene links to create ecosystem-specific functional metagenome

421    predictions.

422    (B-E) For comparison of PICRUSt-KEGG-Ortholog (KO) profiles generated using

423    PICRUST-default and PICRUSt-iMGMC from 16S rRNA gene amplicon sequencing to

424    real KO profiles determined by shotgun metagenome sequencing (WGS) samples

425    from different anatomical locations (n=50) were analyzed.

426    (B) Correlation between KO profiles of metagenomes determined by WGS and

427    PICRUST-default (red) or by WGS and PICRUSt-iMGMC (green) using Pearson and

428    Spearman correlation coefficients. ****: p<0.0001 (two-tailed t-test).

429    (C) Comparison of KO profiles generated using PICRUST-default (red), PICRUSt-

430    iMGMC (green) and WGS (blue) from different anatomical locations. Non-metric

431    multidimensional scaling (NMDS) was performed to visualize similarities.

432    (D) False positive rates and true positive rates were obtained by comparing the

433    PICRUSt-default (red) and PICRUSt-iMGMC (green) KEGG Module predictions

434    against WGS results. The true positive rate reflects the fraction of KEGG Modules

435    commonly predicted by both WGS and PICRUSt default/PICRUSt-iMGMC and the

436    false positive rate reflects the fraction of KEGG Modules that are predicted by PICRUSt

437    default/PICRUSt-IMGMC, but were completely absent in WGS data.

438    (E) KEGG module predictions that differ between PICRUSt-default and PICRUSt-

439    iMGMC predictions. KEGG Module prediction by PICRUSt-default and PICRUSt-

440    iMGMC was compared against WGS for all samples and significant differences in

441    completeness were identified using a Wilcoxon test (FDR-corrected). The heatmap

442    displays select KEGG Modules with highly similar completeness between PICRUSt-

443    iMGMC and WGS, but divergent completeness between PICRUSt-default and WGS
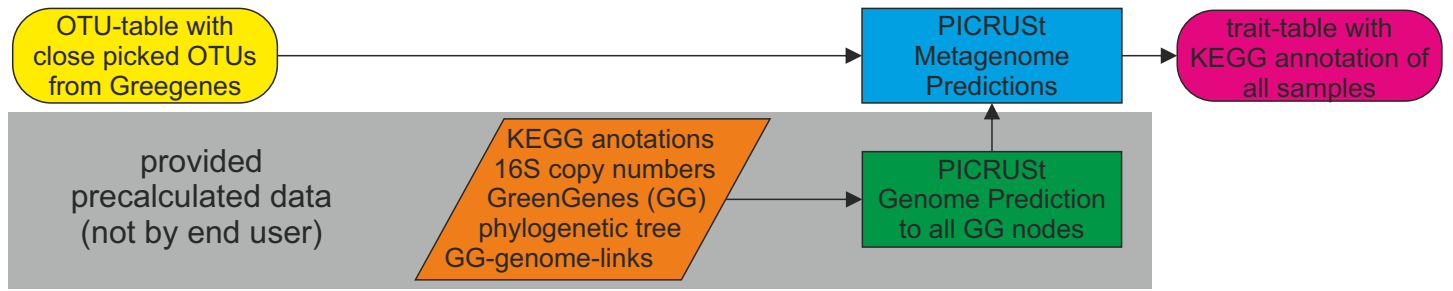
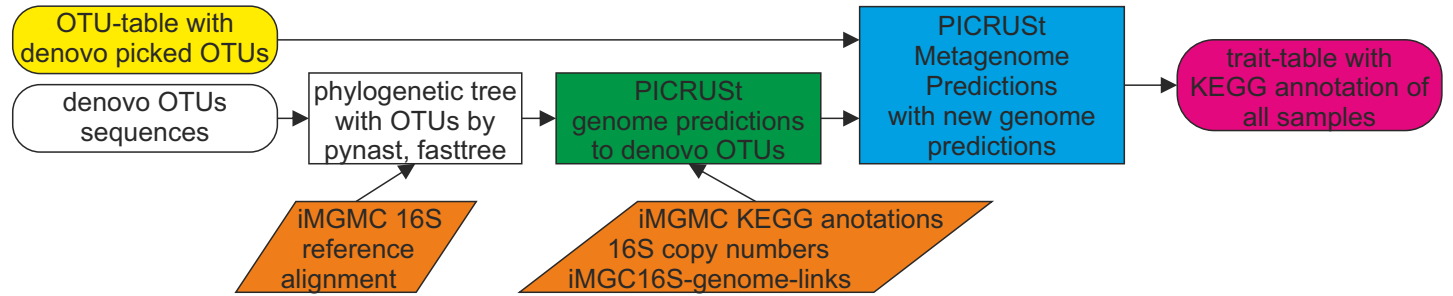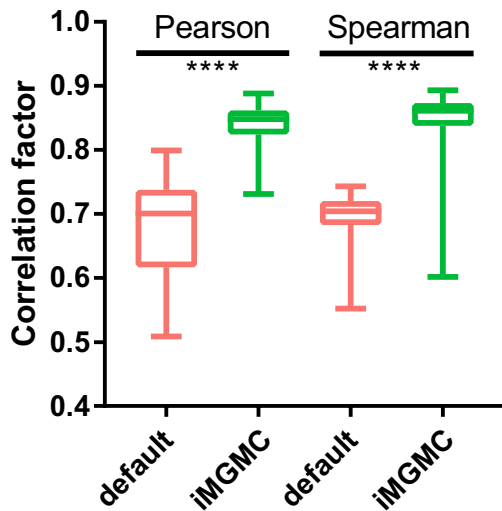444    (see methods for details).

445

446

447

14
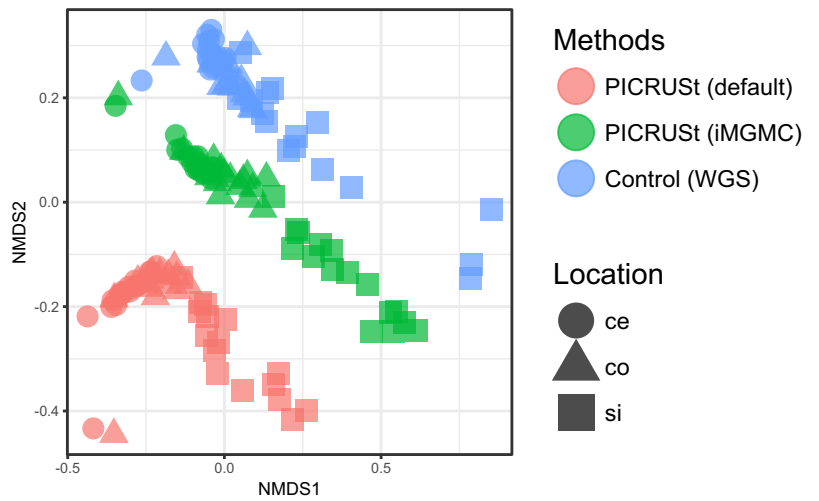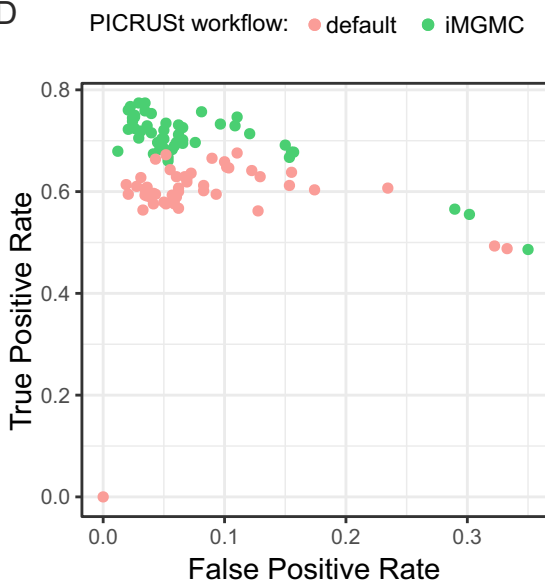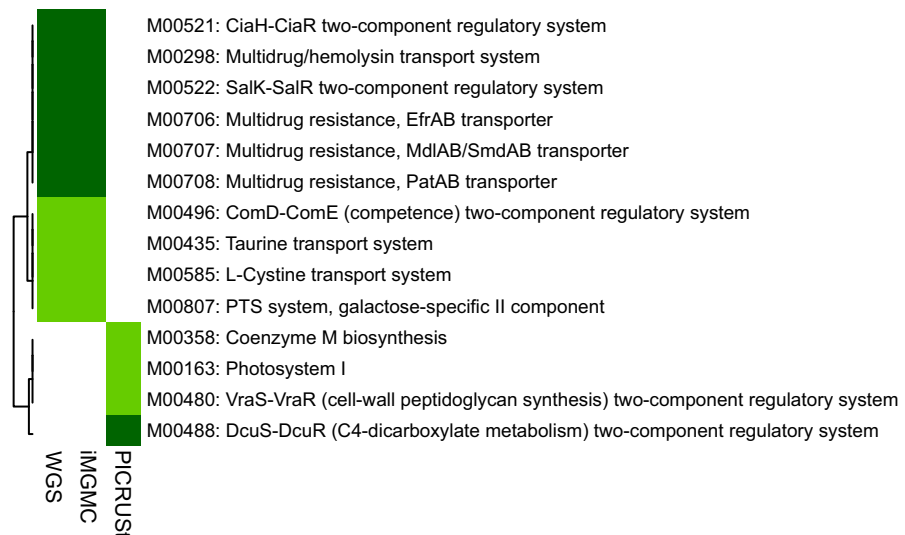
Figure 3

448 **Figure 4:**

449 **Identification of MAGs shared between laboratory mice**

450 (A) Prevalence of MAGs (n=660) in samples from 21 mouse providers. MAGs were

451 considered present in a provider if its relative abundance reached at least 0.1% in one

452 sample of the provider. Numberson the left indicate the fraction (%) and taxonomic

453 grouping  (F: Firmicutes, B: Bacteroides, O: Other phyla) of MAGs with an indicated

454 prevalence (Prev). In the right panel MAGs were ranked by prevalence and dashed

455 lines indicate number of MAGs present in >66%, >50% and >20% of providers,

456 respectively..

457 (B) Comparison of maximal abundance between providers for each MAG (n=22)

458 present in at least 2/3 of providers. For each MAG, the bin number, the highest

459 taxonomic assignment based on the manually curated phylogenetic tree and the

460 provider with the highest abundance is listed. Stars indicate MAGs with matches in

461 NCBI RefSeq.

462 (C) Comparison of the relative abundance of 16S rRNA gene sequences linked to

463 MAGs in the IMNGS database. For each 16S rRNA gene, the closest named relative

464 16S rRNA gene sequence was determined and blasted to the NCBI-16S rRNA gene

465 database. Color of dots and names indicate their taxonomic association to different

466 phyla (F: Firmicutes, B: Bacteroidetes, O: other phyla)

467 (D and E) IMNGS was used to determine the prevalence for iMGMC 16S rRNA gene

468 sequences (n=1,323) in distinct hosts and ecosystems. Of these 1,113 reached at least

469 a prevalence threshold of 1% prevalence within one of the evaluated environment

470 (0.1% sample-depth cutoff of presence). Resulting sequences (n=1,113) were filtered

471 further to have at least >1% relative mean abundance in at least one environment.

472 (D) Heatmap displaying the mean relative abundance within an ecosystem (row

473 normalized) of those 16S rRNA gene sequences which have at least >1% relative

474 mean abundance in at least one environment (n=739).

475 (E) Venn diagram visualizing the distribution of 16S rRNA gene sequences

476 subsampled to be enriched (>50% relative abundance normalized over the

477 ecosystems in Figure 4D) in mouse gut, mouse skin, rat gut and human gut

478 microbiome (n = 569). Numbers indicate fraction of 16S rRNA gene sequences

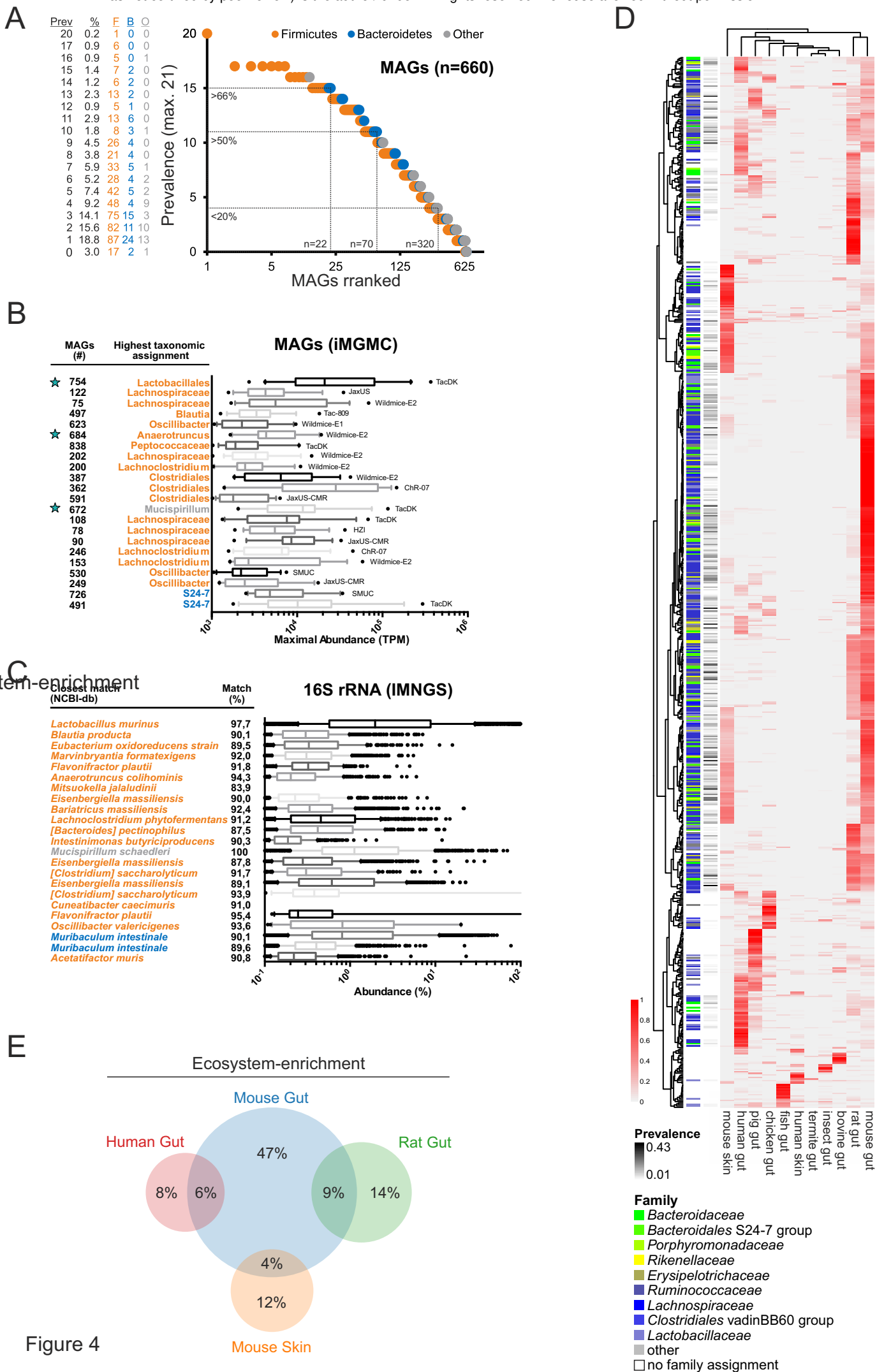479 enriched or shared between indicated ecosystems.

480

Figure 4

**References:**

481 **References:**

482  1.  Kamada, N., Seo, S.-U., Chen, G. Y. & Núñez, G. Role of the gut microbiota in

483      immunity and inflammatory disease. *Nat. Rev. Immunol.* **13,** 321–35 (2013).

484  2.  Li, J. *et al.* An integrated catalog of reference genes in the human gut

485      microbiome. *Nat. Biotechnol.* **32,** 834–841 (2014).

486  3.  Sunagawa, S. *et al.* Structure and function of the global ocean microbiome.

487      *Science (80-. ).* **348,** 1261359–1261359 (2015).

488  4.  Xiao, L. *et al.* A catalog of the mouse gut metagenome. *Nat. Biotechnol.* **33,**

489      1103–1108 (2015).

490  5.  Xiao, L. *et al.* A reference gene catalogue of the pig gut microbiome. *Nat.*

491      *Microbiol.* **1,** 16161 (2016).

492  6.  Lagkouvardos, I. *et al.* The Mouse Intestinal Bacterial Collection (miBC)

493      provides host-specific insight into cultured diversity and functional potential of

494      the gut microbiota. *Nat. Microbiol.* **1,** 16131 (2016).

495  7.  Li, D. *et al.* MEGAHIT v1.0: A fast and scalable metagenome assembler driven

496      by advanced methodologies and community practices. *Methods* **102,** 3–11

497      (2016).

498  8.  Zhu, W., Lomsadze, A. & Borodovsky, M. Ab initio gene identification in

499      metagenomic sequences. *Nucleic Acids Res.* **38,** e132–e132 (2010).

500  9.  Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the

501      next-generation sequencing data. *Bioinformatics* **28,** 3150–3152 (2012).

502  10. Kang, D. D., Froula, J., Egan, R. & Wang, Z. MetaBAT, an efficient tool for

503      accurately reconstructing single genomes from complex microbial

504      communities. *PeerJ* **3,** e1165 (2015).

505  11. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W.

506      CheckM: assessing the quality of microbial genomes recovered from isolates,

507      single cells, and metagenomes. *Genome Res.* **25,** 1043–1055 (2015).

508  12. Miller, C. S., Baker, B. J., Thomas, B. C., Singer, S. W. & Banfield, J. F.

509      EMIRGE: reconstruction of full-length ribosomal genes from microbial

510      community short read sequencing data. *Genome Biol.* **12,** R44 (2011).

511  13. Zeng, F., Wang, Z., Wang, Y., Zhou, J. & Chen, T. Large-scale 16S gene

512      assembly using metagenomics shotgun sequences. *Bioinformatics* **33,** 1447–

513      1456 (2017).

514  14. Delmont, T. O. & Eren, A. M. Identifying contamination with advanced

515      visualization and analysis practices: metagenomic approaches for eukaryotic

516      genome assemblies. *PeerJ* **4,** e1839 (2016).

517  15. Mikheenko, A., Saveliev, V. & Gurevich, A. MetaQUAST: evaluation of

518       metagenome assemblies. *Bioinformatics* **32,** 1088–1090 (2016).

519   16.   Alneberg, J. *et al.* Binning metagenomic contigs by coverage and composition.
520       *Nat. Methods* **11,** 1144–1146 (2014).

521   17.   Parks, D. H. *et al.* Recovery of nearly 8,000 metagenome-assembled
522       genomes substantially expands the tree of life. *Nat. Microbiol.* **2,** 1533–1542
523       (2017).

524   18.   Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome
525       phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36,** 996–1004
526       (2018).

527   19.   Clavel, T., Lagkouvardos, I., Blaut, M. & Stecher, B. The mouse gut
528       microbiome revisited: From complex diversity to model ecosystems. *Int. J.*
529       *Med. Microbiol.* **306,** 316–327 (2016).

530   20.   Karst, S. M. *et al.* Retrieval of a million high-quality, full-length microbial 16S
531       and 18S rRNA gene sequences without primer bias. *Nat. Biotechnol.* **36,** 190–
532       195 (2018).

533   21.   Langille, M. G. I. *et al.* Predictive functional profiling of microbial communities
534       using 16S rRNA marker gene sequences. *Nat. Biotechnol.* **31,** 814–821
535       (2013).

536   22.   Aßhauer, K. P., Wemheuer, B., Daniel, R. & Meinicke, P. Tax4Fun: predicting
537       functional profiles from metagenomic 16S rRNA data: Fig. 1. *Bioinformatics*
538       **31,** 2882–2884 (2015).

539   23.   Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
540       DIAMOND. *Nat. Methods* **12,** 59–60 (2015).

541   24.   Suez, J. *et al.* Artificial sweeteners induce glucose intolerance by altering the
542       gut microbiota. *Nature* **514,** 181–186 (2014).

543   25.   Everard, A. *et al.* Microbiome of prebiotic-treated mice reveals novel targets
544       involved in host response during obesity. *ISME J.* **8,** 2116–2130 (2014).

545   26.   Levy, M. *et al.* Microbiota-Modulated Metabolites Shape the Intestinal
546       Microenvironment by Regulating NLRP6 Inflammasome Signaling. *Cell* **163,**
547       1428–1443 (2015).

548   27.   Sczyrba, A. *et al.* Critical Assessment of Metagenome Interpretation—a
549       benchmark of metagenomics software. *Nat. Methods* **14,** 1063–1071 (2017).

550   28.   Cambuy, Diego D, Coutinho, Felipe H, Dutilh, B. E. Contig annotation tool
551       CAT robustly classifies assembled metagenomic contigs and long sequences.
552       *BioRxiv* 072868 (2016). doi:10.1101/072868

553   29.   Rausch, P. *et al.* Analysis of factors contributing to variation in the C57BL/6J
554       fecal microbiota across German animal facilities. *Int. J. Med. Microbiol.* **306,**

555      (2016).

556   30.   Lagkouvardos, I. *et al.* IMNGS: A comprehensive open resource of processed

557      16S rRNA microbial profiles for ecology and diversity studies. *Sci. Rep.* **6,**

558      33721 (2016).

559