

1 **VIRGO, a comprehensive non-redundant gene catalog, reveals extensive within**
2 **community intraspecies diversity in the human vagina**

3

4 Bing Ma¹, Michael France¹, Jonathan Crabtree¹, Johanna B. Holm¹, Mike Humphrys¹,
5 Rebecca Brotman¹, Jacques Ravel^{1,*}

6

7 ¹Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore,
8 MD 21201, United States

9

10 * corresponding author

11 Jacques Ravel: javel@som.umaryland.edu

12 **Abstract**

13 Background

14 Analysis of metagenomic and metatranscriptomic data is complicated and typically
15 requires extensive computational resources. Leveraging a curated reference database
16 of genes encoded by members of the target microbiome can make these analyses more
17 tractable. Unfortunately, there is no such reference database available for the vaginal
18 microbiome.

19 Results

20 In this study, we assembled a comprehensive human vaginal non-redundant gene
21 catalog (VIRGO) from 264 vaginal metagenomes and 416 genomes of urogenital
22 bacterial isolates. VIRGO includes 0.95 million non-redundant genes compiled from a
23 total of 5.5 million genes belonging to 318 unique bacterial species. We show that
24 VIRGO covers more than 95% of the vaginal bacterial gene content in metagenomes
25 from North American, African, and Chinese women. The gene catalog was extensively
26 functionally annotated from 17 diverse protein databases, and importantly taxonomy
27 was assigned through *in silico* binning of genes derived from metagenomic assemblies.
28 To further enable focused analyses of individual genes and proteins, we also clustered
29 the non-redundant genes into vaginal orthologous groups (VOG). The gene-centric
30 design of VIRGO and VOG provides an easily accessible tool to comprehensively
31 characterize the structure and function of vaginal metagenome and metatranscriptome
32 datasets. To highlight the utility of VIRGO, we analyzed 1,507 additional vaginal
33 metagenomes, uncovering an as of yet undetected high degree of intraspecies diversity
34 within and across vaginal microbiota.

35 Conclusions

36 VIRGO offers a convenient reference database and toolkit that will facilitate a more in-
37 depth understanding of the role of vaginal microorganisms in women's health and
38 reproductive outcomes.

39 **Keywords**

40 vaginal microbiome, metagenome and metatranscriptome reference database, non-
41 redundant gene catalog, intraspecies diversity, gene-centric design, protein family
42 catalog, multi-omics data integration

43 **Background**

44 The microbial communities that inhabit the human body play critical roles in the
45 maintenance of health, and dysfunction of these communities is often associated with
46 disease [1]. Taxonomic profiling of the human microbiome via 16S rRNA gene amplicon
47 sequencing has provided critical insight into the potential role of the microbiota in a wide
48 array of common diseases [2-4]. Yet these data routinely fall short of describing the
49 etiology of such microbiome-associated diseases, such as bacterial vaginosis [5, 6],
50 Crohn's disease [7, 8] or psoriasis [9], among others. This is perhaps because while
51 16S rRNA gene sequencing can provide species-level taxonomic profiles of a microbial
52 community, it does not describe the genes or metabolic functions that are encoded in
53 the constituents' genomes. This is an important distinction because strains of a bacterial
54 species have been documented to exhibit substantial diversity in gene content [10],
55 such that their genomes harbor sets of accessory genes whose presence is variable
56 [11, 12]. It is therefore difficult, if not impossible, to infer the complete function of a
57 microbial species in a specific environment using only the sequence of their 16S rRNA
58 gene. As a consequence, to investigate the role of the human microbiome in health and
59 diseases, particular emphasis should be placed on describing the gene content and
60 gene expression of these microbial communities.

61
62 Metagenomic and metatranscriptomic profiling are emerging approaches aimed at
63 characterizing the gene content and expression of microbial communities. Results have
64 led to increased appreciation for the important role microbial communities play in human
65 health and diseases [13, 14]. Despite the rapid development and increased throughput
66 of sequencing technologies, current knowledge of the genetic and functional diversity of
67 microbial community is still highly limited. This is due, at least in part, to a lack of
68 resources necessary for the analysis of these massive short read datasets [13, 15]. *De*
69 *nov*o assembly of metagenomic or metatranscriptomic datasets typically requires rather
70 substantial computational resources and complicates integration of metagenomic and
71 metatranscriptomic data.

72

73 Accurate, high-resolution mapping of metagenomic or metatranscriptomic data against
74 a comprehensive and curated gene database is an alternative analytical strategy that is
75 less computationally demanding, prone to fewer errors, and provides a standard point of
76 reference for comparison of these data. Development of such curated databases is
77 crucial to further our understanding of the structure and function of microbial
78 communities [15, 16]. In the last two decades, international initiatives such as MetaHit,
79 the NIH funded Human Microbiome Project (HMP) and the International Human
80 Microbiome Consortium (IHMC) were established to generate the resources necessary
81 to enable investigations of the human microbiome, including large reference taxonomic
82 surveys and metagenomic datasets [13, 17]. While multiple 16S rRNA gene catalogs
83 such as RDP [18], SILVA [19], Greengenes [20], and EZBioCloud [21] exist, there are
84 relatively few curated resources for referencing metagenomes and metatranscriptomes.
85 Those that do exist focus only on the gut microbiome of either humans [16, 22] or
86 animal model species [23, 24]. A definite unmet demand exists for reference gene
87 catalogs for other body sites such as the oral cavity, the skin, and the vagina [25].
88
89 In this study, we constructed the human vaginal non-redundant gene catalog (VIRGO),
90 an integrated and comprehensive resource to establish taxonomic and functional
91 profiling of vaginal microbiomes from metagenomic and metatranscriptomic datasets.
92 VIRGO was constructed using 211 *in-house* metagenomes and 53 metagenomes that
93 were generated under the HMP project [26]. The metagenomic data was supplemented
94 with 321 complete or draft genome sequences of urogenital bacterial isolates. The
95 genes identified in the metagenomes and whole genome sequences were further
96 clustered into Vaginal Orthologous Groups (VOGs), a catalog of functional protein
97 families common to vaginal microbiomes. We meticulously curated the gene catalog
98 with taxonomic assignments as well as functional features using 17 diverse protein
99 databases. Importantly, we show that VIRGO provides >95% coverage of the human
100 vaginal microbiome, and it is applicable to populations from North America, Africa and
101 Asia. Together, VIRGO and VOG represent a comprehensive reference repository and
102 a convenient cataloging tool for fast and accurate characterization of vaginal
103 metagenomes and metatranscriptomes. The gene catalog is a compilation of vaginal

104 bacterial species pan-genomes, creating a vaginal “meta-pan-genome”. We further
105 used VIRGO to characterize the amount of intraspecies diversity present in individual
106 vaginal communities. Previous characterization of these communities using either 16S
107 rRNA gene taxonomic profiling or assembly based metagenomic analyses has failed to
108 resolve this diversity. Here we show that vaginal communities contain far more
109 intraspecies diversity than originally expected. This observation challenges the notion
110 that the vaginal microbiota dominated is by one species of *Lactobacillus*, comprised of a
111 single strain, and could have major implications for the ecology of these otherwise low-
112 diversity bacterial communities. Ultimately, VIRGO and its associated analytical
113 framework will facilitate and standardize the analysis and interpretation of large
114 metagenomic and metatranscriptomic datasets thus expanding our understanding of the
115 role of vaginal microbial communities in health and disease.

116 Results

117 VIRGO is sourced from a comprehensive collection of vaginal metagenomes and 118 bacterial genomes

119 VIRGO was constructed using sequence data from fully de-identified vaginal
120 metagenomes (n=264) as well as complete and draft genomes of urogenital bacterial
121 isolates (n=321, de-replicated from 416 genomes). The majority (n=211) of the included
122 metagenomes were sequenced in-house from de-identified vaginal swab specimens. Of
123 the ~18 billion reads generated for these metagenomes, 14.4 billion (79.7%) were
124 identified as human sequences and removed. Interestingly, the proportion of human
125 reads in the vaginal metagenomes was found to vary with community composition.
126 Vaginal metagenomes dominated by *Lactobacillus* spp. had significantly higher
127 proportions of human sequence reads than those from *Lactobacillus* deficient
128 metagenomes (88.7% vs 73.3%; $t=-6.6$, $P < 0.001$; **Additional file 1: Figure S1**).
129 Further pre-processing steps culled sequence reads matching rRNA genes and low
130 sequence quality reads, removing another 1.4% reads. Each metagenome was then *de*
131 *novo* assembled totaling 1.2 million contigs of length > 500bp with a combined length of
132 2.8 billion bp and an N50 of 6.2 kbp. Additional metagenomic data (n=53) were obtained
133 from the HMP [13, 14] and contributed 40,000 contigs with length > 500bp, comprising
134 100 million bp of assembled sequence. The *in-house* metagenomes provided 19.5
135 times more assembled length than the HMP vaginal metagenomes. In addition to the
136 vaginal metagenomes, we also included 321 complete or draft genome sequences of
137 urogenital bacterial isolates, including 139 from HMP and 277 from GenBank and
138 IMG/M (Integrated Microbial Genomes & Microbiomes) [27]. A summary of the
139 metagenomic reads, assembled contigs and genomes included in the construction of
140 VIRGO can be found in **Additional file 2: Table S1**.

141
142 Taxonomic analysis of the 264 metagenomes included in VIRGO, revealed that these
143 communities contained 312 bacterial species present in $\geq 0.01\%$ relative abundance
144 (**Additional file: Table S2**). All major vaginal *Lactobacillus* species (*L. crispatus*, *L.*
145 *gasseri*, *L. iners*, and *L. jensenii*), as well as common facultative and strict anaerobic
146 vaginal species such as *G. vaginalis*, *A. vaginae*, *P. amnii*, *P. timonensis*, *Megasphaera*

147 genomosp., *Mobiluncus mulieris*, *Mageebacillus indolicus* (aka. BVAB3), *Veillonella*
148 *parvula*, among others were identified in the metagenomes. Even BV-associated
149 bacteria that are often only present at low abundance [28] were represented in the
150 metagenomes, including *Finegoldia magna*, *Peptoniphilus harei*, *Peptostreptococcus*
151 *anaerobius*, *Mobiluncus curtisii*, *Peptoniphilus lacrimalis*, *Anaerococcus tetradius*,
152 *Eggerthella* spp., *Ureaplasma urealyticum*, *Veillonella atypica*, *Corynebacterium*
153 *glucuronolyticum*, among others. The taxonomic profiles of these communities were
154 further shown to encompass the five previously reported vaginal community state types
155 (CSTs) [29], CST I, II, III, IV, and V with frequencies in this set of metagenomes of
156 18.9%, 3.8%, 20.5%, 48.5%, and 8.3%, respectively (**Additional file 1: Figure S2.**
157 **Additional File 2: Table S2**). These results highlight the taxonomic breadth of the
158 vaginal bacterial communities included in the construction of VIRGO (**Additional file 1:**
159 **Figure S3**).

160
161 The dataset used to build VIRGO was compiled from vaginal metagenomes that were
162 obtained from North American women. To determine the comprehensiveness of
163 VIRGO, we mapped reads from 91 vaginal metagenomes that were not included in its
164 construction. These metagenomes were obtained from North American, African [30],
165 and Chinese [31] women, allowing us to determine the utility of VIRGO to analyze
166 metagenomes from other populations. Reads from these metagenomes were mapped
167 to the complete and subsets of the sequence contigs used to build VIRGO. More than
168 99% of the reads from North American metagenomes were able to be mapped to the
169 complete VIRGO dataset, while only ~55% of these reads mapped to contigs from the
170 HMP vaginal metagenomes subset (**Fig. 1, Additional file 2: Table S3**). This result
171 indicates a lack of genetic diversity in the HMP vaginal metagenomes, which were
172 derived from highly selected and healthy women [32]. Further, despite originating from
173 populations not used in the construction of VIRGO, 96% and 88% of the reads from
174 African and Chinese women mapped to the complete VIRGO dataset. For these two
175 cohorts, 71.7% and 99.9% of the reads that failed to map to VIRGO, also did not have a
176 match in GenBank (**Additional file 1: Figure S4**). These results illustrate the
177 comprehensiveness of VIRGO and its broad application to different populations and

178 ethnicities. It further shows that the bacterial genetic diversity in the vaginal microbiome
179 across populations is somewhat homogenous.

180

181 VIRGO: a non-redundant vaginal bacterial gene catalog

182 Coding sequences (CDS, n=5,509,298) were predicted from the metagenomic
183 assemblies and genome sequences using MetageneMark [33]. The core workflow to
184 identify and cluster these CDSs is shown in **Fig. 2**, and a more detailed illustration is
185 provided in **Additional file 1: Figure S5**. Metagenomic assemblies contributed ~80% of
186 the CDSs while the remaining ~20% of CDSs originated from the urogenital bacteria
187 isolate genome sequences. Redundant genes were then identified and removed via a
188 greedy pairwise comparison at the nucleotide level using highly stringent criteria of 95%
189 identity over 90% of the shorter gene length [16, 22]. This process afforded the removal
190 of partial genes and eliminated overcalling genes as unique because of sequencing
191 errors. A total of 948,158 non-redundant CDSs longer than 99 bp were identified and
192 retained, representing 17.2% of the original 5.5 million CDSs. The *in-house* vaginal
193 metagenomes used to build VIRGO contributed 12 times more non-redundant genes
194 (634,288 genes) than the HMP vaginal metagenomes (54,500 genes). Combined, the
195 metagenomes contributed twice as many non-redundant genes as urogenital bacterial
196 isolate genome sequences (371,099 genes). Metagenomes were found to contain a
197 higher proportion of redundant genes than bacterial genome sequences (84.5% versus
198 58.1% of their sequence lengths) (**Additional file 2: Table S3**).

199

200 In order to facilitate the use of VIRGO to characterize vaginal microbial communities,
201 each non-redundant gene was taxonomically and functionally annotated. Non-
202 redundant genes were assigned to taxonomic groups using a custom pipeline as
203 depicted in **Additional file 1: Figure S5**. First, metagenomic contigs were assigned
204 taxonomy if 95% of the composite reads were annotated to the same species. Second,
205 genes encoded on an metagenomic contig with assigned taxonomy were given the
206 taxonomy of that contig (details in Methods). A total of 458,526 non-redundant genes
207 comprising 48.4% of VIRGO were able to be taxonomically curated. Overall, 269 unique
208 bacterial species were annotated in VIRGO (**Additional file 2: Table S4**), representing

209 a majority of the described vaginal species (**Additional file 1: Figure S2**). This includes
210 BVAB1, an as of yet unculturable vaginal species, for which several metagenome-
211 assembled genomes (MAGs) were recently made available (accession # will be
212 provided upon acceptance of the manuscript). BVAB1 was only been previously
213 detectable using a partial 16S rRNA gene reference sequence [34]. It was found
214 abundantly present in most of the metagenomes with a prevalence of 15.6% and mean
215 abundance of (18.9% +/- 0.01) as shown in **Figure 3D**. When stratified by CST, CST IV
216 metagenomes have the smallest proportion (<30%) of their gene content taxonomically
217 annotated (**Additional file 1: Figure S6**) compared to ~45-50% in *Lactobacillus*-
218 dominated CSTs. The most abundant species based on gene content are shown in **Fig.**
219 **3a** and **Additional file 1: Figure S7**. Besides bacteria, we also curated potential fungal
220 and phage genes (details in Methods) that were generally present in low abundance if
221 detected at $0.17 \pm 0.04\%$ and $0.03 \pm 0.001\%$, respectively. An additional 10,908 fungal
222 and 15,965 phage genes were included (**Additional file 2: Table S5**,
223 <https://github.com/Ravel-Laboratory/VIRGO>).

224
225 By including many metagenomes and bacterial isolate genome sequences, we sought
226 to capture each vaginal species' pangenome in VIRGO. To determine the extent to
227 which we were successful, we generated metagenome accumulation curves for the
228 number of non-redundant genes belonging to several key vaginal species (**Fig. 3b**)
229 These curves track the number of new non-redundant genes added when increasing
230 numbers of metagenomes containing a given species are included in constructing the
231 database. The accumulation curves for six of the seven species tested (*L. crispatus*, *L.*
232 *iners*, *L. gasseri*, *L. jensenii*, *P. timonensis*, *A. vaginae*) have reached saturation (**Fig.**
233 **3b**). This indicates that VIRGO includes the majority of these species pangenomes. The
234 number of non-redundant genes included for five out of these six species are similar
235 (~5,000 genes), while the sixth, *A. vaginae*, had twice as many. This pales in
236 comparison to the number of non-redundant genes included in VIRGO for *G. vaginalis*,
237 which surpasses 25,000 genes. *G. vaginalis* is the only species analyzed for which
238 saturation as estimated by metagenome accumulation curves, was not reached.
239

240 The non-redundant genes were decorated with a rich set of functional annotations. We
241 performed intensive functional annotation using both the JCVI standard operating
242 procedure [HMP] for annotating prokaryotic metagenomic shotgun sequencing data [35]
243 as well as 17 additional functional protein databases including KEGG, COG, eggNOG,
244 gene product, CDD, and GO, among others. A complete list of the functional annotation
245 sources employed to characterize the VIRGO non-redundant genes is illustrated in **Fig.**
246 **2**, and an overview of the eggNOG functions encoded in VIRGO is shown in **Fig. 3c**.
247 Overall 785,268 genes (82.8% of all non-redundant genes) were assigned a functional
248 annotation from at least one source. This gene-rich annotation of the non-redundant
249 gene catalog enables a comprehensive functional characterization of vaginal
250 metagenomes and metatranscriptomes.

251 VOG: orthologous protein families in vaginal microbiome

252 The non-redundant genes were translated into amino acid sequences and clustered into
253 vaginal orthologous groups (VOGs). The resulting database of VOGs can be used to
254 interrogate the protein families found in the vaginal microbiome. A modified Jaccard
255 index was used as a measure of similarity between amino acid sequences [36, 37].
256 Briefly, the similarity between each pair of proteins was calculated as the intersection
257 divided by the union of the list of proteins connected to the pair of proteins, (**Fig. 2** and
258 **Additional file 1: Figure S5** algorithm accessible at [https://github.com/Ravel-](https://github.com/Ravel-Laboratory/VIRGO)
259 [Laboratory/VIRGO](https://github.com/Ravel-Laboratory/VIRGO)). The resulting connected graph of proteins is referred as Jaccard
260 clusters (JACs), and reciprocal best hits of JACs is referred as Jaccard orthologous
261 clusters (JOCs) (details in Methods). The JOCs orthologous protein families can be
262 highly conserved (alignment score >950) or partially aligned with both conserved and
263 variable regions (alignment score ~300) (**Additional file 1: Figure S8**). This highlights
264 the flexibility of the network-based aggregation algorithm used to recruit both highly
265 similar and distantly related proteins without imposing a single similarity threshold. A
266 total of 617,127 JACs and 552,679 JOCs were generated, of which 177,684 contained
267 at least two genes while the remaining 374,995 are singletons, indicating 38.5% of all
268 VOG proteins are unique. The sequences, alignment, and phylogenetic trees for each of
269 the JOCs will be available at <https://github.com/Ravel-Laboratory/VIRGO>.

270

271 Complementary to the VIRGO non-redundant gene sequences, VOG provides an amino
272 acid sequence reference that can be used to improve functional annotation,
273 comparative genomics and evolution of vaginal orthologous protein families. For
274 example, we used VOG and retrieved 32 proteins of the orthologous family encoding
275 vaginolysin, a *G. vaginalis* cholesterol-dependent cytolysin that is key to its
276 pathogenicity as it forms pore in epithelial cells [38, 39] (**Additional file 2: Table S6**,
277 **Additional file 1: Figure S9**). Using the retrieved alignment, we identified 3 amino acid
278 variants in a 11-amino acid sequences of domain 4 of vaginolysin, one of the three
279 variants, an alanine-to-valine substitution that is divergent across *G. vaginalis* and had
280 not been reported previously. This example illustrates how VOG can be mined to
281 understand biological relevance and to generate hypotheses. In this case it points to
282 potential differences in pore formation activity and possibly cytotoxicity, which could be
283 further investigated. As another example to use VOG for a large-scale data mining of
284 protein family of interest, we searched VOG using the key phrase “cell surface-
285 associated proteins” and “*L. iners*” and retrieved two protein families, one of which was
286 recognized to have an LPXTG motif while the other harbored the motif YSIRK
287 (**Additional file 2: Table S7**). Interestingly, a previous study on staphylococcal proteins
288 suggested that the motifs LPXTG and YSIRK were involved in different biological
289 processes related to surface protein anchoring to cell wall envelope [40], and both are
290 implicated in virulence by promoting bacterial attachment to alpha- and beta-chains of
291 human fibrinogen and formation of bacterial clumps [41]. These two retrieved protein
292 families are specific to *L. iners* and provide relevant evidence for future experimental
293 validation to understand adherence and related biological processes. These two
294 examples demonstrate how the VOG database can be used to explore more
295 mechanistic understandings of vaginal bacterial communities.

296

297 Gene richness is characteristic of vaginal microbiomes

298 Gene richness, calculated as number of non-redundant genes, has been adapted as the
299 proxy of genetic diversity based on community gene content, and more recently, as
300 community-level biomarker in gut quantitative metagenomics studies [42, 43]. We
301 applied this paradigm to vaginal metagenomes included in VIRGO and defined high

302 gene count (HGC) vaginal communities as those that contained >10,000 non-redundant
303 genes and low gene count (LGC) vaginal communities as those that contained $\leq 10,000$
304 non-redundant genes. The number of non-redundant genes identified in a metagenome
305 was not found to correlate with the depth sequencing (**Figure 3E. Additional file 2:**
306 **Table S8**). As expected, HGC communities had a significantly higher number of non-
307 redundant genes ($29,898 \pm 1,025$) than LGC communities ($4,920 \pm 151.6$), however these
308 types of communities also showed differences in their functional makeup. The LGC
309 communities were found to be enriched for genes related to carbohydrate transport and
310 metabolism, as well as those involved in transcription, while HGC communities were
311 found to be enriched in genes related to intracellular trafficking, secretion, and vesicular
312 transport, including coenzyme transport and metabolism (**Additional file 1: Figure**
313 **S10**). We also found that *Lactobacillus*-dominated communities were typically
314 categorized as LGC (82.9%) and *Lactobacillus*-deficient communities as HGC (88.3%)
315 (**Fig. 4a**). However, this was not always the case, most notably, *L. iners*-dominated
316 communities were classified as HGC 21.7% of the time, the highest percentage among
317 all *Lactobacillus*-dominated communities. In fact, *L. iners*-dominated communities
318 ($7,803 \pm 6,973$) generally had a greater gene richness than *L. crispatus*-dominated
319 ($5,409 \pm 3,392$), *L. gasseri*-dominated ($3,909 \pm 2,761$), and *L. jensenii*-dominated
320 ($3,990 \pm 3,230$) communities. Further, *L. iners* in HGC communities and *L. iners* in LGC
321 communities show distinct functional makeup (**Additional file 1: Figure S11**). Similarly,
322 not all *Lactobacillus*-deficient communities were classified as HGC—11.7% of these
323 communities were identified as LGC. This includes communities with a high abundance
324 of *G. vaginalis*, whose gene richness varied between $7,689 \pm 1,700$ in LGC and
325 $16,887 \pm 566$ in HGC communities.

326

327 In addition to being a characteristic of individual communities, gene richness can also
328 be used to characterize individual genes based on their observed preference for either
329 HGC or LGC communities. Using data from the 264 vaginal metagenomes, we
330 classified each non-redundant gene as either an HGC or LGC gene if $\geq 95\%$ its
331 occurrences were in HGC or LGC communities, respectively. Genes that did not meet
332 this criterion were annotated as having no preference. These gene richness annotations

333 were included for each non-redundant gene in VIRGO. For example, 84.1%, 53.3%,
334 60.5% of top prevalent tryptophan biosynthesis genes in VIRGO, tryptophanase
335 (TNA), tryptophan synthase beta chain (TRPB), and tryptophanyl-tRNA synthetase
336 (TRPS), are HGC genes, while 0%, 0%, and 7.0% are LGC genes (**Additional file1:**
337 **Table S9**). Given the top most affiliated taxonomic groups for these tryptophan
338 biosynthesis genes were identified as *G. vaginalis*, *A. vaginae*, *M. mulieris* (**Additional**
339 **file 1: Figure S12**) our result indicates tryptophan biosynthesis genes are most
340 prevalent in BV-associated bacteria of high gene richness vaginal communities,
341 agreeing with recent studies [44, 45].

342

343 Using these gene annotations, we were further able to evaluate whether a vaginal
344 bacterial species' genes were overrepresented as being HGC or LGC (**Fig. 4b**).
345 *Lactobacillus* spp., particularly *L. crispatus*, *L. jensenii*, *L. gasseri*, *L. vaginalis*, were
346 observed to be highly overrepresented in LGC communities. On the other hand, genes
347 belonging to many other BV-associated species, specifically *P. timonensis*, *P. buccalis*,
348 *P. amnii*, *M. mulieris*, BVAB3, *Porphyromonas uenonis*, *P. harej*, *Anaerococcus*
349 *tetradus*, *M. curtisii*, were overrepresented in HGC. These results demonstrate gene
350 richness category information, characteristics of vaginal metagenomic communities as
351 well as individual genes in the community, provides additional dimension to facilitate our
352 understanding of the genetic basis of the biological processes that drive vaginal
353 microbiomes.

354 Integration of metagenome and metatranscriptome data using VIRGO as a reference 355 framework

356 By serving as a reference, VIRGO enables the characterization and integrative analyses
357 of the abundance of genes and their expression in the vaginal microenvironment. To
358 demonstrate its use, we analyzed a woman's vaginal metagenomes and associated
359 metatranscriptomes at four time points over an episode of symptomatic bacterial
360 vaginosis (BV): prior to (T1), during (T2 & T3), and after (T4) (**Fig. 5a**). Not surprisingly,
361 the expressed functions represented in the metatranscriptomes were often different
362 from the encoded functional makeup of the corresponding metagenomes (**Fig. 5b**). For
363 example, T4 genes related to translation were underrepresented in the

364 metatranscriptome as compared to the metagenome, while genes of unknown function
365 were overrepresented. VIRGO enables rapid binning of genes by species, which
366 revealed dramatic differences in gene abundance and their transcriptional activity in
367 vaginal species (**Fig. 5c**). Prior to the BV episode (T1), a small proportion of *L. iners*
368 genes were present (1.5%) but these genes exhibited high expression levels,
369 accounting for over 20% of the metatranscriptome. At the same time point, *L. crispatus*
370 genes made up the majority of the gene present (96.3%) but exhibited low expression
371 levels (34.2%). In contrast, at the end of the BV episode, *L. crispatus* gene made up a
372 small proportion of the metagenome (T3) but were highly transcriptionally active. This
373 increased activity corresponded with *L. crispatus* regaining dominance at T4, following
374 the resolution of the BV episode. Similarly, despite its low abundance, *P. harei* was
375 highly transcriptionally active during the BV episode (T3), expressing transcript
376 associated with amino acid transport and metabolism, indicating a potential role for this
377 bacterial species in the etiology or symptomology associated with BV. Interestingly, the
378 functional makeup of *G. vaginalis* is similar at T2 and T3, but its metatranscriptome is
379 enriched for functions involved in energy production and conversion at T2, and enriched
380 for functions related to translation, energy production, and carbohydrate metabolism at
381 T3. These examples highlight how VIRGO can be used to integrate metagenome and
382 metatranscriptomic datasets to gain better functional insights into the vaginal
383 microbiome.

384 VIRGO revealed high within-community intraspecies diversity

385 VIRGO can be used to characterize the genome content of individual bacterial species
386 that are present in the vaginal microbiome. We applied VIRGO to a dataset of 1,507 *in-*
387 *house* and publicly available vaginal metagenomes, to characterize the gene content of
388 four *Lactobacillus* species (*L. crispatus*, *L. iners*, *L. jensenii*, and *L. gasseri*) and three
389 additional species commonly found in the vagina (*G. vaginalis*, *A. vaginae* and *P.*
390 *timonensis*). We recovered most of each species gene content (>80% of the average
391 gene count in a genome) even when that species was present at low abundance (<1%)
392 in a community. For instance, even though *P. timonensis* [46] was generally present in
393 low abundance in these metagenomes (4.8% \pm 0.3% mean \pm S.E., range [0.1-33.8%]),
394 we recovered the majority of its genome (2,469 \pm 401 CDS, **Additional file 1: Figure**

395 **S13; Additional file 2: Table S10**). We observed similarly high sensitivity in the
396 analysis of the other six selected vaginal species (**Fig. 6a, Additional file 2: Table**
397 **S10**). These results demonstrate VIRGO's capability for characterizing the gene content
398 of low abundance taxa from metagenomic data.

399
400 Using these species-specific gene repertoires, we characterize the amount of
401 intraspecies diversity present within an individual woman's vaginal microbiome.
402 Because VIRGO basically comprises the "pangenomes" of each vaginal bacterial
403 species, it can be used to evaluate the amount of intraspecies diversity present in these
404 communities. For this analysis, we counted the number of genes that were assigned to
405 each of the seven species in each of the 1,507 metagenomic datasets and compared
406 this number to that found in each species' reference genomes. The number of genes for
407 a species in a community often exceeded that found in a single isolate genome (**Fig. 6a,**
408 **6b**), suggesting that multiple strains of a species co-occur in vaginal bacterial
409 communities. The total number of *L. crispatus* genes identified in each of the
410 metagenomes where it was detected contained on average 1.6 times more genes
411 ($3,262 \pm 586$) than that found encoded on *L. crispatus* genomes ($2,064 \pm 225$, $P < 0.001$).
412 Similar results were observed for *G. vaginalis*, *A. vaginae*, *L. iners*, *L. jensenii*, and *L.*
413 *gasseri*, which are represented by 7.0, 3.4, 2.2, 1.3, and 1.1 times more genes in
414 metagenomes than that found in genomes, respectively. *G. vaginalis* and *A. vaginae*
415 exhibited the highest degree of intraspecies diversity, while *L. crispatus* has the highest
416 within-metagenome intraspecies diversity among all major vaginal *Lactobacillus* spp.
417 (**Additional file 1: Figure S13; Figure 6c**). These results suggest that a woman's
418 vaginal bacterial populations are routinely comprised of more than one strain of most
419 species. VIRGO affords investigating this unprecedented intraspecies diversity in
420 vaginal communities.

421
422 We next applied well-established practices from pangenomics [10, 12] in order to
423 identify core and accessory non-redundant genes among our sample-specific species
424 gene repertoires. Based on the clustering patterns of gene prevalence profiles, we were
425 able to define groups of consistently present (core) and variably present (accessory)

426 non-redundant genes. The majority of the observed genes for each of the species were
427 categorized as accessory, with variable representation across the metagenomic
428 datasets. Using *L. crispatus* as an example, we observed more than twice as many non-
429 redundant genes with variable representation across the metagenomes than those
430 present in every sample (**Fig. 6c**). Interestingly, it is clear from this analysis that the
431 gene content identified with VIRGO in genome sequences of *L. crispatus* under-
432 represent the intraspecies genetic diversity present in the metagenomes. Similar results
433 were observed for the other six species analyzed, although the magnitude of the
434 difference between the metagenome and isolate gene repertoires varied depending on
435 the species. Overall, VIRGO revealed that metagenomic data carry a more extensive
436 gene content than is found in all combined isolate genome sequences.

437 Metagenomic subspecies in vaginal ecosystem

438 Hierarchical clustering of the metagenome species-specific gene content profiles
439 revealed distinct groupings which we term “metagenomic subspecies” (MG-subspecies).
440 These metagenomic subspecies represent types of bacterial populations that share a
441 similar gene pool as assessed by shotgun metagenomic sequence data. For example,
442 this analysis revealed at least three distinct metagenomic subspecies for *L. gasseri*
443 (**Fig. 6d**). *L. gasseri* MG-subspecies I and III have large sets of non-redundant genes
444 that are present in one but not the others, while *L. gasseri* MG-subspecies II carries a
445 blend of the genes from both MG-subspecies I and III. The analysis of *G. vaginalis*
446 revealed more than four types of profile groupings, though concordant with the
447 previously described multiple types of isolate genomes [47], we find that this genome-
448 based paradigm largely under-represents the diversity of *G. vaginalis* gene content
449 identified in metagenomes (**Additional file 1: Figure S13e**). We applied this analysis to
450 seven vaginal species (**Additional file 1: Figure S13**) and found that vaginal microbial
451 communities are often composed of complex mixtures of multiple strains of the same
452 species, and that these mixtures can be clustered into distinct MG-subspecies. Further
453 interrogation of these vaginal MG-subspecies and their gene content is likely to reveal
454 novel features of vaginal communities and their sub-populations that will contribute to
455 our understanding of the vaginal ecosystem of niche-optimized strains.

456 Discussion

457 Microbiome studies have become increasingly sophisticated with the rapid
458 advancement of sequencing throughput and the associated decrease in sequencing
459 cost. However, identifying features that drive correlations between the microbiome and
460 health using multi-omics sequence data remains challenging. This is due, in part, to
461 difficulties in analyzing and integrating the complex, feature rich, metagenomic and
462 metatranscriptomic data now common to microbiome studies. A scalable tool that
463 provides a comprehensive characterization of such multi-omics data is therefore highly
464 desired. VIRGO is a large vaginal microbiome database designed to fulfill such research
465 needs for investigations of the vaginal microbiome and its relation to women's health. In
466 summary, VIRGO has (i) a comprehensive breadth that includes previously observed
467 community types, vaginal species, and even fungi and viruses; (ii) a gene-centric design
468 that enables the integration of functional and taxonomic characterization of
469 metagenomic and metatranscriptomic data originating from the same sample; (iii) a high
470 scalability and low memory requirement; (iv) a high sensitivity that affords
471 characterization of the gene content of low-abundance bacteria; (v) an easy to use
472 framework from which to evaluate gene richness and within-species diversity.

473
474 VIRGO contains a multitude of non-redundant genes that we identified in vaginal
475 metagenomes and urogenital bacterial isolates. These non-redundant genes were also
476 clustered into orthologous groups (VOGs) using a memory-efficient network-based
477 algorithm that handles nodes connectivity in high dimensionality space [48, 49]. This
478 approach to identifying orthologous protein sequences allows for great flexibility
479 because it does not rely on a single sequence similarity cutoff value [50, 51]. These
480 families of vaginal orthologs will assist the development of a mechanistic understanding
481 of these proteins and how they relate to health. For example, van der Veer and co-
482 workers recently identified and characterized the *L. crispatus* pullulanase (*pulA*) gene
483 which they show encodes an enzyme with amylase activity that likely allows this species
484 to degrade host glycogen in the vaginal environment [52]. Using VIRGO and VOG, we
485 were able to identify pullulanase domain containing proteins in 37 other vaginal taxa
486 including: *G. vaginalis*, *L. iners* and *P. timonensis* (**Additional file 2: Table S12**),

487 providing insight into the breadth of vaginal bacteria that may be capable of degrading
488 host glycogen. In this way, VIRGO and VOG can facilitate knowledge retrieval,
489 hypothesis generation and future experimental validation to advance understanding of
490 vaginal ecosystem.

491
492 Using VIRGO, we observed that vaginal metagenomes varied in gene richness, with
493 some communities having more non-redundant genes than others. Gene richness has
494 been found to be indicative of the pathophysiological state of the gut microbiome in
495 studies of obesity [43], dietary intervention [42], type II diabetes [53], and inflammation
496 and metabolic disease [54]. We adapted the concept of gene richness as a
497 characterization of community gene content and defined an analogous definition for the
498 vaginal microbiome. An outstanding difference in gene richness was observed between
499 *Lactobacillus*-dominated and *Lactobacillus*-deficient communities. Approximately 85%
500 of communities with a high relative abundance of *Lactobacillus* sp., had a low gene
501 richness across the community, whereas *Lactobacillus*-deficient communities were
502 more likely to have a high gene richness. However, around 22% of *Lactobacillus*-
503 dominated communities did have high gene richness and 12% of *Lactobacillus*-deficient
504 communities had low gene richness. It may be that gene richness category, when
505 combined with community state types, provides a useful and, ecologically relevant,
506 categorization of vaginal community states. For example, it is envisioned that a subject
507 that has a *Lactobacillus*-dominated community with high gene richness is at a higher
508 risk of switching to a dysbiotic state than one whose community is dominated by
509 *Lactobacillus* but with low gene richness. In such case, VIRGO provides the analytical
510 suite needed to test this and other hypotheses relating gene richness to the ecology of
511 the vaginal microbiome.

512
513 In our demonstrative analysis of more than 1,500 metagenomes, we identified and
514 characterized a wealth of intraspecies diversity that was present within individual
515 vaginal microbial communities. Populations of bacterial species in vaginal communities
516 comprises of multiple strains. Previous studies of the vaginal microbiome have largely
517 treated these species as singular genotypes [55, 56], although some more recent

518 studies have examined intraspecies diversity in these communities [57, 58].
519 Intraspecies diversity is important because it is likely to influence many properties of the
520 communities including their temporal stability and resilience, as well as how they relate
521 to host health. Unfortunately, intraspecies diversity is difficult to detect using typical
522 assembly-based metagenomic analysis strategies, which are notoriously ill suited for
523 resolving strains of the same species [59, 60]. VIRGO can be a more suitable tool for
524 characterizing intraspecies diversity because it was built to contain the non-redundant
525 pangenomes of most species common to the vagina. Strict mapping of sequence reads
526 against the VIRGO database provides an accurate and sensitive way of identifying the
527 aggregated non-redundant genes that belongs to each species in a metagenome. We
528 expect VIRGO to facilitate future investigations of intraspecies diversity in vaginal
529 microbial communities. We further showed that, for the seven species we examined, the
530 intraspecies diversity had structure. Vaginal metagenomes from different subjects
531 contained related sets of species-specific non-redundant genes. We postulate that
532 these clusters of samples with shared gene content represent similar collectives of
533 strains which we have termed “metagenomic subspecies”. It is expected that, given their
534 shared gene content, these metagenomic subspecies might also share phenotypic
535 characteristics. However, additional studies are needed to characterize differences
536 between metagenomic subspecies and to detail their possible effect on host health.
537 Reconstructing a particular species’ metagenomic subspecies might be possible by
538 identifying and combining isolates that adequately cover the genetic repertoire of the
539 metagenomic subspecies. One complication to this approach is that, in many cases, the
540 observed metagenomic subspecies contain non-redundant genes which have not been
541 observed in isolate genome sequences for that species. This could reflect a limitation in
542 the number of isolate genomes available for a species or even systematic bias in the
543 growth and recovery of species from the vagina [61]. Targeted isolation of strains from
544 communities containing the metagenomic subspecies of interest are needed in order to
545 fill in these gaps in the future.
546
547 The value of VIRGO resides in its functions as both a central repository and a highly
548 scalable tool for fast, accurate characterization of vaginal microbiomes. VIRGO is

549 particularly useful for users with limited computational skills, a large volume of
550 sequencing data, and/or limited computing infrastructure. In particular, the
551 metagenome-metatrascriptome data integration enabled by the gene-centric design in
552 VIRGO provides a powerful approach to determine the expression patterns of microbial
553 functions, and in doing so, to characterize contextualized complex mechanisms of host-
554 microbiota interactions in vaginal communities. This feature makes possible the meta-
555 analyses of vaginal microbiome features and the quantitative integration of findings from
556 multiple studies, which helps with the common issue of confounding gene copy number
557 that has been a major challenge in analyzing metatranscriptomic dataset [62, 63]. We
558 also anticipate that VIRGO will be used to process metaproteomic datasets when that
559 practice becomes common and easily accessible. Each of the protein sequence of each
560 gene could be used to map peptides obtained from metaproteomic pipelines and access
561 VIRGO rich annotation. On the other hand, we acknowledge the limitations of the
562 referenced based approach of VIRGO. This version is focused on the gene-level de-
563 redundancy and characterization of vaginal microbiome. However, in the future we plan
564 to expand VIRGO to include the capability to identify nucleotide variants within a gene.
565 We believe this will further facilitate our understanding of within-species diversity and
566 evolutionary change in the vaginal ecosystem. The database is primarily focused on
567 bacteria with limited inclusion of viral and fungal gene sequences. Future in-depth
568 profiling of these non-bacterial microbes will allow VIRGO to provide a more complete
569 picture of vaginal microbial communities.

570

571

572

573 **Conclusion**

574 Efforts are underway to translate our growing understanding of human-associated
575 microbial communities into clinical biomarkers and treatments. A deeper understanding
576 of the complex mechanisms of host-microbiota interactions requires the integration of
577 multi-omics data. VIRGO presents a central reference database and analytical
578 framework to enable the efficient and accurate characterization of the microbial gene
579 content of the human vaginal microbiome. Powered by a rich suite of functional and
580 taxonomic annotations, VIRGO allows for the integrated analysis of metagenomic and
581 metatranscriptomic data. VIRGO further provides a gene-centric approach to describe
582 vaginal microbial community structure including fine scale variation at the intraspecies
583 level. This unprecedented view of intraspecies diversity within a vaginal community is
584 far beyond the scope offered by current genome references. VIRGO is a centralized,
585 and freely available resource for vaginal microbiome studies. It will facilitate the analysis
586 of multi-omics data now common to microbiome studies, and provide comprehensive
587 insight into community membership, function, and ecological perspective of the vaginal
588 microbiome.
589

590 **Methods**

591 Datasets

592 Metagenomes used in this study include 211 newly *in-house* sequenced datasets and
593 53 vaginal datasets downloaded from the HMP data repository ([http://www.hmpdacc-](http://www.hmpdacc-resources.org/cgi-bin/hmp_catalog/main.cgi)
594 [resources.org/cgi-bin/hmp_catalog/main.cgi](http://www.hmpdacc-resources.org/cgi-bin/hmp_catalog/main.cgi)). Genome sequences of urogenital
595 bacterial isolates deposited in multiple databases were downloaded on November 10,
596 2016, including GenBank (<http://www.ncbi.nlm.nih.gov/>), IMG/M: Integrated Microbial
597 Genomes & Microbiomes (<https://img.jgi.doe.gov/>), and HMP referencing genome
598 database (http://www.hmpdacc-resources.org/cgi-bin/hmp_catalog/main.cgi). After
599 removing duplicate genomes under the same strain names, genomes of 416 urogenital
600 bacterial strains and 321 bacterial species were included in the catalog. A full list of the
601 genomes and metagenomes used in the construction of the database can be found in
602 **Additional file: Table S1.**

603 Nucleic acid extraction, library construction, and metagenome and metatranscriptome 604 sequencing.

605 The included 211 *in-house* metagenomes were generated as follows: whole genomic
606 DNA was extracted from 300 µl aliquot of vaginal ESwab re-suspended into 1ml Amies
607 transport medium (ESwab, Copan Diagnostics Inc.) and preserved at -80°C. Briefly,
608 Cells were then lysed using a combination of enzymatic digestion and mechanical
609 disruption that included mutanolysin, lysostaphin and lysozyme treatment, followed by
610 proteinase K, SDS and bead beating steps. Procedures for DNA extraction and
611 concentration qualification were previously described [29, 64]. The shotgun
612 metagenomic sequence libraries were constructed from the extracted DNA using
613 Illumina Nextera XT kits and sequenced on an Illumina HiSeq 2500 platform at the
614 Genomic Resource Center at the University of Maryland School of Medicine.

615

616 The metatranscriptomes used to demonstrate the use of VIRGO for the analysis of
617 community-wide gene expression were obtained from RNA extracted from vaginal
618 swabs stored in 2 ml Amies Transport Medium-RNA later solution (50%/50%, vol/vol)
619 archived at -80°C. A total of 500 µl of ice-cold PBS was added to 1,000 µl of that
620 solution and spun down at 8,000xg for 10 min. The pellet was resuspended in 500 µl

621 ice-cold RNase-free PBS with 10 μ l β -mercaptoethanol. The suspension was
622 transferred to Lysis Matrix B tube (MP Biomedicals) containing 100 μ l 10% SDS and
623 500 μ l acid phenol and beads beaded using a FastPrep instrument (MP Biomedicals)
624 for 45 seconds at 5.5 m/s. The aqueous phase was mixed with 250 μ l acid phenol and
625 250 μ l 24:1 chloroform:isoamyl alcohol. The aqueous layer was again transferred to a
626 fresh tube and mixed with 500 μ l 24:1 chloroform:isoamyl alcohol. For every 300 μ l
627 resulting aqueous solution, we added 30 μ l of 3 M sodium acetate, 3 μ l of glycogen (5
628 mg/ml), and three volumes of 100% ethanol. The mixture was incubated at -20°C
629 overnight to precipitate the nucleic acids. After centrifugation at 13,400xg for 30 min at
630 4°C, the resulting pellet was washed, dried, and dissolved in 100 μ l of DEPC-treated
631 water. Carryover DNA was removed by: 1) treating twice with Turbo DNase free
632 (Ambion, Cat. No. AM1907) at two half-hour intervals according to the manufacturer's
633 protocol for rigorous DNase treatment, 2) purifying twice using gDNA-eliminator
634 columns (QIAGEN) before and after DNase treatment followed by RNeasy column
635 purification (QIAGEN). We further conducted PCR using 16S rRNA primer 27F (5'-
636 AGAGTTTGATCCTGGCTCAG -3') and 534R (5'- CATTACCGCGGCTGCTGG -3') to
637 confirm DNA removal. The quality of extracted RNA was checked using an Agilent 2100
638 Expert Bioanalyzer Nano chip. Ribosomal RNA removal was performed according to the
639 manufacturer's protocol of a combined Gram-positive, Gram-negative and
640 Human/mouse/rat Ribo-Zero rRNA Removal Kit (Epicentre Technologies). The resulting
641 RNA was purified using Zymo RNA clean & Concentrator-5 column kit (ZYMO
642 Research). RNA final quality was checked using an Agilent RNA 6000 Expert
643 Bioanalyzer Pico chip. Sequencing libraries of A and B containing 6 bp indexes were
644 prepared using the TruSeq RNA sample prep kit (Illumina) following a modification of
645 the manufacturer's protocol: cDNA was purified between enzymatic reactions and
646 library size selection was performed with AMPure XT beads (Beckman Coulter
647 Genomics). Library sequencing was performed using the Illumina HiSeq 2500 platform.

648 Construction of the human vaginal non-redundant gene catalog (VIRGO)

649 Multiple bioinformatics pre-processing steps were applied to the raw shotgun
650 metagenomic sequence datasets, including (1) eliminating all human sequence reads
651 (including human rRNA LSU/SSU sequence reads) using BMTagger v3.101 [65] against

652 a standard human genome reference (GRCh37.p5 [66]), (2) *in silico* microbial rRNA
653 sequence reads depletion by aligning all reads using Bowtie (v1) [67] against the SILVA
654 PARC ribosomal-subunit sequence database [19] to eliminate mis-assemblies of these
655 repeated regions. After each of these steps, the paired reads were removed; (3)
656 stringent quality control using Trimmomatic [68], in which the Illumina adapter was
657 trimmed and reads with average quality greater than Q15 using a sliding window of 4 bp
658 with no ambiguous base calling were retained. MetaPhlAn (v2) [69] was subsequently
659 used to establish taxonomic profiles after these pre-processing steps. Samples were
660 then clustered in community state types (CSTs) using taxa abundance tables and the
661 Jensen-Shannon divergence metrics as previously described [29, 70]. Species
662 accumulation curves and diversity estimates for rarefied samples were computed using
663 R package *iNEXT* [71] and *vegan* [72]. The 264 vaginal metagenomes were then
664 assembled using IDBA-UD (v1.0) [73] with a k value range of 20-100. Genes were
665 called on the resulting contigs using MetageneMark (v3.25) [33] to predict CDSs with
666 the default settings. Genes and gene fragments that were at least 99bp long, with
667 greater than 95% identity over 90% of the shorter gene length were clustered together
668 by a greedy pairwise comparison implemented in CD-HIT-EST (v4.6) [74], according to
669 the clustering procedure and threshold defined previously [16, 22]. The gene with the
670 longest length ≥ 99 bp was used as the representative for each cluster of redundant
671 genes.

672 Taxonomic and functional annotations of VIRGO

673 The non-redundant genes were annotated with a rich set of taxonomic and functional
674 information. Genes that originated from an isolate sequence genome were automatically
675 assigned that species name. For metagenomes, taxonomy was assigned to a
676 metagenomic contig by mapping the sequence reads making up that contig to the
677 Integrated Microbial Genomes (IMG) reference database (v400) using bowtie (v1,
678 parameters: “-l 25 --fullref --chunkmbs 512 --best --strata -m 20”). A secondary filter was
679 applied so that the total number of mismatches between the read and the reference was
680 less than 35, and that the first 25 bp of the read matched the reference. Using the
681 results of this mapping, taxonomy was assigned to all genes encoded on the contig that
682 met the following four criteria: 1) at least 95% of the reads mapped to the same

683 bacterial species, 2) the remaining 5% off-target reads did not map to a single species,
684 3) the contig had at least 2X average coverage and >50 reads, 4) at least 25% of the
685 contig length had reads mapped onto. These stringent criteria were used to ensure high
686 fidelity of the taxonomic assignments and a low contribution of potentially chimeric
687 contigs. To further diminish the risk of incorporating false taxonomic assignments, the
688 annotations of the contigs belonging to species at low relative abundance in the sample
689 were removed. Genome completeness was estimated as the fractional representation of
690 the genome in the metagenome using BLASTN (minimal overlapping >60% of the
691 shorter sequence and >80% sequence similarity). For each metagenome, only
692 taxonomic assignments originating from species with at least 80% representation were
693 incorporated. The genes that shows >80% sequence similarity over 60% of query gene
694 length to the non-redundant genes were then assigned. The non-redundant genes in
695 VIRGO were searched against fungal database that includes 5 vaginal yeast species in
696 40 genomes (listed in **Additional file 2: Table S5**) using BLASTN, that a gene must
697 have at least 80% sequencing similarity with over 60% overlapping length to be curated.
698 We also annotated potential phage genes that may be present in VIRGO by searching
699 against phage orthologous groups or Prokaryotic virus orthologous groups (version
700 2016) [51, 75], using BLASTN and included the ones at >80% sequence similarity over
701 60% of query gene length in annotation (**Additional file 2: Table S5**). Functional
702 annotations based on the standard procedure for each of 17 functional databases,
703 including: cluster of orthologous groups (COG[76], eggNOG (v4.5) [77], KEGG[78]),
704 conserved protein domain (CDD[79], Pfam[80], ProDom[81], PROSITE[82],
705 TIGRFAM[83], InterPro[84]), domain architectures (CATH-Gene3D[85, 86],
706 SMART[87]), intrinsic protein disorder (MobiDB[88]), high-quality manual annotation
707 (HAMAP[89]), protein superfamily (PIRSF[90]), a compendium of protein fingerprints
708 (PRINTS[91]), and gene product attributes (Gene Ontology [92], JCVI SOP [35]).

709

710 Construction of vaginal orthologous groups (VOGs) for protein families

711 The non-redundant genes were also clustered based on orthology to generate a set of
712 Vaginal Orthologous Groups (VOGs). To do this we used a modified version of a
713 Jaccard clustering method previously implemented [36, 37]. We performed an all-

714 versus-all BLASTP search among the translated coding sequences (CDS) of the non-
715 redundant genes included in VIRGO [93, 94]. The all-against-all BLASTP matches was
716 used to compute Jaccard similarity coefficient for each pair of translated CDSs, without
717 constraints based on which sample or microorganism from which it originated. Only
718 BLASTP matches with 80% sequence identity and 70% overlap, and an E-value less
719 than 1E-10 were used in the calculation of the Jaccard similarity coefficient. The filtered
720 BLASTP results were then used to define connections between pairs of translated
721 CDSs resulting in a network graph with the translated CDSs as nodes and their
722 connections as edges. The Jaccard similarity coefficient was then calculated as the
723 number of nodes that had direct connections to the two translated CDSs divided by the
724 total number of nodes that had direct connections to either of the two translated CDSs
725 in the network (intersection and union) [37]. Jaccard clusters (JACs) were defined as a
726 set of translated CDSs whose Jaccard similarity coefficient was at least 0.55. If two
727 translated CDSs from different JACs were reciprocal best matches according to the
728 BLASTP searches, the two JACs were merged. Finally, the alignment program T-Coffee
729 [95] was used to assess the alignment quality within the JACs and to calculate the
730 alignment score.

731 Bioinformatics analysis

732 The comprehensiveness of VIRGO was tested using vaginal metagenomic data from
733 vaginal metagenomes of North American women not including in the construction of
734 VIRGO and sequenced in this study, as well as women from African [30], and China
735 [31]. The sequences reads were first mapped to the VIRGO contigs using bowtie (v2;
736 parameters: --threads 4 --sensitive-local -D 10 -R 2 -N 0 -L 22 -i S,1,1.75 -k 1 --ignore-
737 quals --no-unal) [96], according to the criteria used previously in the construction of a
738 gut gene catalog [16]). Any unmapped reads were compared to the GenBank nt
739 database [97] using BLASTN and an E-value of 1E-10 as cutoff. To annotate BVAB1
740 genes in VIRGO, we used BLASTN and an E-value of 1E-10 as cutoff, the matched
741 genes with percent identity >95% over >90% of gene length were annotated as BVAB1
742 genes. To retrieve pullulanase (*pulA*) genes in VIRGO, we used conserved protein
743 domain CDD [79] annotation and keyword “pullulanase”. To further demonstrate the
744 comprehensiveness of VIRGO and that VIRGO captures the pangenome of selected

745 species, species specific metagenome accumulation curves for the number of non-
746 redundant genes were constructed for seven vaginal species by rarefaction with 100
747 bootstraps: *L. crispatus*, *L. iners*, *L. jensenii*, *L. gasseri*, and *G. vaginalis*, *A. vaginae*
748 and *P. timonensis*.

749

750 For gene count category and analysis, the included 264 vaginal metagenomes were
751 classified as either having a high gene count (>10,000 non-redundant genes) or low
752 gene count (<10,000 non-redundant genes). The VIRGO non-redundant genes were
753 then annotated as either being a high gene count gene or low gene count gene if the
754 gene was preferentially identified (at least 95%) in high or low gene count
755 metagenomes. The log ratio of genes of a species being in either high or low gene
756 count metagenomes across the 264 vaginal metagenomes was calculated for all
757 species with at least 0.1% abundance and at least 100 genes in either HGC or LGC
758 groups. The species with more than 4 times more abundant (in logarithm 2 scale) in a
759 category (either HGC or LGC) were considered showing preference in one of the
760 categories.

761 Using VIRGO to characterize within community intraspecies diversity

762 Intraspecies diversity analyses were conducted by mapping isolate genome sequences
763 as well as vaginal metagenomes to VIRGO. We chose to focus our analysis on the
764 previously mentioned seven vaginal species. Accession numbers for genomes of the
765 four *Lactobacillus* species (*L. crispatus*, *L. iners*, *L. jensenii*, and *L. gasseri*) and three
766 additional species (*G. vaginalis*, *A. vaginae* and *P. timonensis*) can be found in
767 **Additional file 2: Table S11**. A total of 1,507 vaginal metagenomes including 1,403 *in-*
768 *house* from de-identified vaginal swab and lavage specimens and 76 publicly available,
769 were mapped against VIRGO. Their accession numbers can be found in **Additional file**
770 **2: Table S12**. For each of the seven species, a presence/absence matrix for the
771 species' non-redundant genes was constructed that included the data from species'
772 isolate genomes and all metagenomes that contained at least 80% of the average
773 number of genes encoded on a genome of that species. Comparisons of the number of
774 non-redundant genes present in the species isolate genomes versus the metagenomes
775 in which they appeared were conducted using student t-test. Hierarchical clustering

776 was performed on the boolean matrix of the species' non-redundant genes using
777 Jaccard clustering implemented in the *vegan* package in R [98]. A tutorial describing
778 how to use VIRGO and VOG is available online at [https://github.com/Ravel-](https://github.com/Ravel-Laboratory/VIRGO)
779 [Laboratory/VIRGO](https://github.com/Ravel-Laboratory/VIRGO).

780 **List of abbreviations**

- 781 **VIRGO**: human vaginal non-redundant gene catalog
782 **VOG**: vaginal orthologous groups
783 **IMG/M**: Integrated Microbial Genomes & Microbiomes
784 **rRNA**: ribosomal ribonucleic acid
785 **CST**: vaginal community state type
786 **MAG**: metagenome-assembled genome
787 **JAC**: Jaccard cluster
788 **JOC**: Jaccard orthologous cluster
789 **HGC**: high gene count
790 **LGC**: low gene count
791 **MG-subspecies**: metagenomic subspecies
792 **Av**: *A. vaginae*
793 **Gv**: *G. vaginalis*
794 **Pt**: *P. timonensis*
795 **Lc**: *L. crispatus*
796 **Li**: *L. iners*
797 **Lj**: *L. jensenii*
798 **Lg**: *L. gasseri*
799

800 References

- 801 1. Cho I, Blaser MJ: **The human microbiome: at the interface of health and**
802 **disease.** *Nat Rev Genet* 2012, **13**:260-270.
- 803 2. Henao-Mejia J, Elinav E, Thaïss CA, Licona-Limon P, Flavell RA: **Role of the**
804 **intestinal microbiome in liver disease.** *J Autoimmun* 2013, **46**:66-73.
- 805 3. Ley RE: **Obesity and the human microbiome.** *Curr Opin Gastroenterol* 2010,
806 **26**:5-11.
- 807 4. Zeeuwen PL, Kleerebezem M, Timmerman HM, Schalkwijk J: **Microbiome and**
808 **skin diseases.** *Curr Opin Allergy Clin Immunol* 2013, **13**:514-520.
- 809 5. Nasioudis D, Linhares IM, Ledger WJ, Witkin SS: **Bacterial vaginosis: a critical**
810 **analysis of current knowledge.** *BJOG* 2017, **124**:61-69.
- 811 6. Schwebke JR: **New concepts in the etiology of bacterial vaginosis.** *Curr*
812 *Infect Dis Rep* 2009, **11**:143-147.
- 813 7. Gevers D, Kugathasan S, Knights D, Kostic AD, Knight R, Xavier RJ: **A**
814 **Microbiome Foundation for the Study of Crohn's Disease.** *Cell Host Microbe*
815 2017, **21**:301-304.
- 816 8. Wright EK, Kamm MA, Teo SM, Inouye M, Wagner J, Kirkwood CD: **Recent**
817 **advances in characterizing the gastrointestinal microbiome in Crohn's**
818 **disease: a systematic review.** *Inflamm Bowel Dis* 2015, **21**:1219-1228.
- 819 9. Tett A, Pasolli E, Farina S, Truong DT, Asnicar F, Zolfo M, Beghini F, Armanini F,
820 Jousson O, De Sanctis V, et al: **Unexplored diversity and strain-level**
821 **structure of the skin microbiome associated with psoriasis.** *NPJ Biofilms*
822 *Microbiomes* 2017, **3**:14.
- 823 10. Tettelin H, Massignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli
824 SV, Crabtree J, Jones AL, Durkin AS, et al: **Genome analysis of multiple**
825 **pathogenic isolates of Streptococcus agalactiae: implications for the**
826 **microbial "pan-genome".** *Proc Natl Acad Sci U S A* 2005, **102**:13950-13955.
- 827 11. Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R: **The microbial pan-**
828 **genome.** *Current opinion in genetics & development* 2005, **15**:589-594.
- 829 12. Tettelin H, Riley D, Cattuto C, Medini D: **Comparative genomics: the bacterial**
830 **pan-genome.** *Current opinion in microbiology* 2008, **11**:472-477.
- 831 13. Consortium TH: **A framework for human microbiome research.** *Nature* 2012,
832 **486**:215-221.
- 833 14. Consortium THMP: **Structure, function and diversity of the healthy human**
834 **microbiome.** *Nature* 2012, **486**:207-214.
- 835 15. Human Microbiome Jumpstart Reference Strains C, Nelson KE, Weinstock GM,
836 Highlander SK, Worley KC, Creasy HH, Wortman JR, Rusch DB, Mitreva M,
837 Sodergren E, et al: **A catalog of reference genomes from the human**
838 **microbiome.** *Science* 2010, **328**:994-999.
- 839 16. Qin J, Li R, Raes J, Arumugam M, Burgdorf KS, Manichanh C, Nielsen T, Pons
840 N, Levenez F, Yamada T, et al: **A human gut microbial gene catalogue**
841 **established by metagenomic sequencing.** *Nature* 2010, **464**:59-65.
- 842 17. Group NHW, Peterson J, Garges S, Giovanni M, McInnes P, Wang L, Schloss
843 JA, Bonazzi V, McEwen JE, Wetterstrand KA, et al: **The NIH Human**
844 **Microbiome Project.** *Genome Res* 2009, **19**:2317-2323.

- 845 18. Cole JR, Chai B, Farris RJ, Wang Q, Kulam-Syed-Mohideen AS, McGarrell DM,
846 Bandela AM, Cardenas E, Garrity GM, Tiedje JM: **The ribosomal database**
847 **project (RDP-II): introducing myRDP space and quality controlled public**
848 **data.** *Nucleic Acids Res* 2007, **35**:D169-172.
- 849 19. Quast C, Priesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J,
850 Glockner FO: **The SILVA ribosomal RNA gene database project: improved**
851 **data processing and web-based tools.** *Nucleic acids research* 2013, **41**:D590-
852 596.
- 853 20. McDonald D, Price MN, Goodrich J, Nawrocki EP, DeSantis TZ, Probst A,
854 Andersen GL, Knight R, Hugenholtz P: **An improved Greengenes taxonomy**
855 **with explicit ranks for ecological and evolutionary analyses of bacteria and**
856 **archaea.** *The ISME journal* 2012, **6**:610-618.
- 857 21. Yoon SH, Ha SM, Kwon S, Lim J, Kim Y, Seo H, Chun J: **Introducing**
858 **EzBioCloud: a taxonomically united database of 16S rRNA gene sequences**
859 **and whole-genome assemblies.** *Int J Syst Evol Microbiol* 2017, **67**:1613-1617.
- 860 22. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, Arumugam M, Kultima JR,
861 Prifti E, Nielsen T, et al: **An integrated catalog of reference genes in the**
862 **human gut microbiome.** *Nat Biotechnol* 2014, **32**:834-841.
- 863 23. Xiao L, Feng Q, Liang S, Sonne SB, Xia Z, Qiu X, Li X, Long H, Zhang J, Zhang
864 D, et al: **A catalog of the mouse gut metagenome.** *Nat Biotechnol* 2015,
865 **33**:1103-1108.
- 866 24. Xiao L, Estelle J, Kiillerich P, Ramayo-Caldas Y, Xia Z, Feng Q, Liang S,
867 Pedersen AO, Kjeldsen NJ, Liu C, et al: **A reference gene catalogue of the pig**
868 **gut microbiome.** *Nat Microbiol* 2016:16161.
- 869 25. Wu H, Tremaroli V, Backhed F: **Linking Microbiota to Human Diseases: A**
870 **Systems Biology Perspective.** *Trends Endocrinol Metab* 2015, **26**:758-770.
- 871 26. Turnbaugh PJ, Ley RE, Hamady M, Fraser-Liggett CM, Knight R, Gordon JI: **The**
872 **human microbiome project.** *Nature* 2007, **449**:804-810.
- 873 27. Ribeiro FJ, Przybylski D, Yin S, Sharpe T, Gnerre S, Abouelleil A, Berlin AM,
874 Montmayeur A, Shea TP, Walker BJ, et al: **Finished bacterial genomes from**
875 **shotgun sequence data.** *Genome Res* 2012, **22**:2270-2277.
- 876 28. Fredricks DN, Fiedler TL, Marrazzo JM: **Molecular identification of bacteria**
877 **associated with bacterial vaginosis.** *N Engl J Med* 2005, **353**:1899-1911.
- 878 29. Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SS, McCulle SL, Karlebach S,
879 Gorle R, Russell J, Tacket CO, et al: **Vaginal microbiome of reproductive-age**
880 **women.** *Proc Natl Acad Sci USA* 2011, **108 Suppl 1**:4680-4687.
- 881 30. Gosmann C, Anahtar MN, Handley SA, Farcasanu M, Abu-Ali G, Bowman BA,
882 Padavattan N, Desai C, Droit L, Moodley A, et al: **Lactobacillus-Deficient**
883 **Cervicovaginal Bacterial Communities Are Associated with Increased HIV**
884 **Acquisition in Young South African Women.** *Immunity* 2017, **46**:29-37.
- 885 31. Chen C, Li F, Wei W, Wang Z, Dai J, Hao L, Song L, Zhang X, Zeng L, Du H, et
886 al: **The metagenome of the female upper reproductive tract.** *Gigascience*
887 2018.
- 888 32. Gevers D, Knight R, Petrosino JF, Huang K, McGuire AL, Birren BW, Nelson KE,
889 White O, Methe BA, Huttenhower C: **The Human Microbiome Project: a**

- 890 **community resource for the healthy human microbiome. *PLoS Biol* 2012,**
891 **10:e1001377.**
- 892 33. Zhu W, Lomsadze A, Borodovsky M: **Ab initio gene identification in**
893 **metagenomic sequences. *Nucleic acids research* 2010, 38:e132.**
- 894 34. Fredricks DN, Fiedler TL, Marrazzo JM: **Molecular identification of bacteria**
895 **associated with bacterial vaginosis. *N Engl J Med* 2005, 353:1899-1911.**
- 896 35. Tanenbaum DM, Goll J, Murphy S, Kumar P, Zafar N, Thiagarajan M, Madupu R,
897 Davidsen T, Kagan L, Kravitz S, et al: **The JCVI standard operating procedure**
898 **for annotating prokaryotic metagenomic shotgun sequencing data.**
899 ***Standards in genomic sciences* 2010, 2:229-237.**
- 900 36. Riley DR, Angiuoli SV, Crabtree J, Dunning Hotopp JC, Tettelin H: **Using Sybil**
901 **for interactive comparative genomics of microbes on the web.**
902 ***Bioinformatics* 2012, 28:160-166.**
- 903 37. Crabtree J, Angiuoli SV, Wortman JR, White OR: **Sybil: methods and software**
904 **for multiple genome comparison and visualization. *Methods in molecular***
905 ***biology* 2007, 408:93-108.**
- 906 38. Tweten RK: **Cholesterol-dependent cytolysins, a family of versatile pore-**
907 **forming toxins. *Infect Immun* 2005, 73:6199-6209.**
- 908 39. Gelber SE, Aguilar JL, Lewis KL, Ratner AJ: **Functional and phylogenetic**
909 **characterization of Vaginolysin, the human-specific cytolysin from**
910 ***Gardnerella vaginalis*. *J Bacteriol* 2008, 190:3896-3903.**
- 911 40. Bae T, Schneewind O: **The YSIRK-G/S motif of staphylococcal protein A and**
912 **its role in efficiency of signal peptide processing. *J Bacteriol* 2003,**
913 **185:2910-2919.**
- 914 41. Cogen AL, Nizet V, Gallo RL: **Skin microbiota: a source of disease or**
915 **defence? *Br J Dermatol* 2008, 158:442-455.**
- 916 42. Cotillard A, Kennedy SP, Kong LC, Prifti E, Pons N, Le Chatelier E, Almeida M,
917 Quinquis B, Levenez F, Galleron N, et al: **Dietary intervention impact on gut**
918 **microbial gene richness. *Nature* 2013, 500:585-588.**
- 919 43. Le Chatelier E, Nielsen T, Qin J, Prifti E, Hildebrand F, Falony G, Almeida M,
920 Arumugam M, Batto JM, Kennedy S, et al: **Richness of human gut microbiome**
921 **correlates with metabolic markers. *Nature* 2013, 500:541-546.**
- 922 44. Ziklo N, Vidgen ME, Taing K, Huston WM, Timms P: **Dysbiosis of the Vaginal**
923 **Microbiota and Higher Vaginal Kynurenine/Tryptophan Ratio Reveals an**
924 **Association with *Chlamydia trachomatis* Genital Infections. *Front Cell Infect***
925 ***Microbiol* 2018, 8:1.**
- 926 45. Aiyar A, Quayle AJ, Buckner LR, Sherchand SP, Chang TL, Zea AH, Martin DH,
927 Belland RJ: **Influence of the tryptophan-indole-IFN γ axis on human**
928 **genital *Chlamydia trachomatis* infection: role of vaginal co-infections. *Front***
929 ***Cell Infect Microbiol* 2014, 4:72.**
- 930 46. Beamer MA, Austin MN, Avolia HA, Meyn LA, Bunge KE, Hillier SL: **Bacterial**
931 **species colonizing the vagina of healthy women are not associated with**
932 **race. *Anaerobe* 2017, 45:40-43.**
- 933 47. Ahmed A, Earl J, Retchless A, Hillier SL, Rabe LK, Cherpes TL, Powell E, Janto
934 B, Eutsey R, Hiller NL, et al: **Comparative genomic analyses of 17 clinical**
935 **isolates of *Gardnerella vaginalis* provide evidence of multiple genetically**

- 936 **isolated clades consistent with subspeciation into genovars. *Journal of*
937 *bacteriology* 2012, **194**:3922-3937.**
- 938 48. Mirkin BG, Fenner TI, Galperin MY, Koonin EV: **Algorithms for computing**
939 **parsimonious evolutionary scenarios for genome evolution, the last**
940 **universal common ancestor and dominance of horizontal gene transfer in**
941 **the evolution of prokaryotes.** *Bmc Evolutionary Biology* 2003, **3**:-
- 942 49. Fletcher MN, Castro MA, Wang X, de Santiago I, O'Reilly M, Chin SF, Rueda
943 OM, Caldas C, Ponder BA, Markowitz F, Meyer KB: **Master regulators of**
944 **FGFR2 signalling and breast cancer risk.** *Nat Commun* 2013, **4**:2464.
- 945 50. Kristensen DM, Kannan L, Coleman MK, Wolf YI, Sorokin A, Koonin EV,
946 Mushegian A: **A low-polynomial algorithm for assembling clusters of**
947 **orthologous groups from intergenomic symmetric best matches.**
948 *Bioinformatics* 2010, **26**:1481-1487.
- 949 51. Kristensen DM, Waller AS, Yamada T, Bork P, Mushegian AR, Koonin EV:
950 **Orthologous gene clusters and taxon signature genes for viruses of**
951 **prokaryotes.** *J Bacteriol* 2013, **195**:941-950.
- 952 52. van der Veer C, Hertzberger R, Bruisten S, Tytgat H, Swanenburg J, de Kat
953 Angelino-Bart A, Schuren F, Molenaar D, Reid G, de Vries H, Kort R:
954 **Comparative genomics of human Lactobacillus crispatus isolates reveals**
955 **genes for glycosylation and glycogen degradation: Implications for in vivo**
956 **dominance of the vaginal microbiota.** *bioRxiv* 2018.
- 957 53. Karlsson FH, Tremaroli V, Nookaew I, Bergstrom G, Behre CJ, Fagerberg B,
958 Nielsen J, Backhed F: **Gut metagenome in European women with normal,**
959 **impaired and diabetic glucose control.** *Nature* 2013, **498**:99-103.
- 960 54. Boulange CL, Neves AL, Chilloux J, Nicholson JK, Dumas ME: **Impact of the**
961 **gut microbiota on inflammation, obesity, and metabolic disease.** *Genome*
962 *Med* 2016, **8**:42.
- 963 55. Martin DH, Marrazzo JM: **The Vaginal Microbiome: Current Understanding**
964 **and Future Directions.** *J Infect Dis* 2016, **214 Suppl 1**:S36-41.
- 965 56. Greenbaum S, Greenbaum G, Moran-Gilad J, Weintraub AY: **Ecological**
966 **dynamics of the vaginal microbiome in relation to health and disease.** *Am J*
967 *Obstet Gynecol* 2018.
- 968 57. Goltsman DSA, Sun CL, Proctor DM, DiGiulio DB, Robaczewska A, Thomas BC,
969 Shaw GM, Stevenson DK, Holmes SP, Banfield JF, Relman DA: **Metagenomic**
970 **analysis with strain-level resolution reveals fine-scale variation in the**
971 **human pregnancy microbiome.** *Genome Res* 2018, **28**:1467-1480.
- 972 58. Deng ZL, Gottschick C, Bhuju S, Masur C, Abels C, Wagner-Dobler I:
973 **Metatranscriptome Analysis of the Vaginal Microbiota Reveals Potential**
974 **Mechanisms for Protection against Metronidazole in Bacterial Vaginosis.**
975 *mSphere* 2018, **3**.
- 976 59. Albanese D, Donati C: **Strain profiling and epidemiology of bacterial species**
977 **from metagenomic sequencing.** *Nat Commun* 2017, **8**:2260.
- 978 60. Segata N: **On the Road to Strain-Resolved Comparative Metagenomics.**
979 *mSystems* 2018, **3**.
- 980 61. Sommer MO: **Advancing gut microbiome research using cultivation.** *Curr*
981 *Opin Microbiol* 2015, **27**:127-132.

- 982 62. Franzosa EA, Hsu T, Sirota-Madi A, Shafquat A, Abu-Ali G, Morgan XC,
983 Huttenhower C: **Sequencing and beyond: integrating molecular 'omics' for**
984 **microbial community profiling.** *Nat Rev Microbiol* 2015, **13**:360-372.
- 985 63. Giannoukos G, Ciulla DM, Huang K, Haas BJ, Izard J, Levin JZ, Livny J, Earl
986 AM, Gevers D, Ward DV, et al: **Efficient and robust RNA-seq process for**
987 **cultured bacteria and complex community transcriptomes.** *Genome Biol*
988 2012, **13**:R23.
- 989 64. Yuan S, Cohen DB, Ravel J, Abdo Z, Forney LJ: **Evaluation of methods for the**
990 **extraction and purification of DNA from the human microbiome.** *PLoS One*
991 2012, **7**:e33865.
- 992 65. Rotmistrovsky K, Agarwala R: **BMTagger: Best Match Tagger for removing**
993 **human reads from metagenomics datasets.** NCBI/NLM, National Institutes of
994 Health; 2011.
- 995 66. Church DM, Schneider VA, Graves T, Auger K, Cunningham F, Bouk N, Chen
996 HC, Agarwala R, McLaren WM, Ritchie GR, et al: **Modernizing reference**
997 **genome assemblies.** *PLoS biology* 2011, **9**:e1001091.
- 998 67. Langmead B, Trapnell C, Pop M, Salzberg SL: **Ultrafast and memory-efficient**
999 **alignment of short DNA sequences to the human genome.** *Genome biology*
1000 2009, **10**:R25.
- 1001 68. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina**
1002 **sequence data.** *Bioinformatics* 2014, **30**:2114-2120.
- 1003 69. Segata N, Waldron L, Ballarini A, Narasimhan V, Jousson O, Huttenhower C:
1004 **Metagenomic microbial community profiling using unique clade-specific**
1005 **marker genes.** *Nature methods* 2012, **9**:811-814.
- 1006 70. Gajer P, Brotman RM, Bai G, Sakamoto J, Schutte UM, Zhong X, Koenig SS, Fu
1007 L, Ma ZS, Zhou X, et al: **Temporal dynamics of the human vaginal**
1008 **microbiota.** *Sci Transl Med* 2012, **4**:132ra152.
- 1009 71. Chao A, Ma KH, Hsieh TC: **iNEXT (iNterpolation and EXTrapolation) Online:**
1010 **Software for Interpolation and Extrapolation of Species Diversity.**; 2016.
- 1011 72. Dixon P: **VEGAN, a package of R functions for community ecology.** *Journal*
1012 *of Vegetation Science* 2003, **14**:927-930.
- 1013 73. Peng Y, Leung HC, Yiu SM, Chin FY: **IDBA-UD: a de novo assembler for**
1014 **single-cell and metagenomic sequencing data with highly uneven depth.**
1015 *Bioinformatics* 2012, **28**:1420-1428.
- 1016 74. Li W, Jaroszewski L, Godzik A: **Clustering of highly homologous sequences**
1017 **to reduce the size of large protein databases.** *Bioinformatics* 2001, **17**:282-
1018 283.
- 1019 75. Graziotin AL, Koonin EV, Kristensen DM: **Prokaryotic Virus Orthologous**
1020 **Groups (pVOGs): a resource for comparative genomics and protein family**
1021 **annotation.** *Nucleic Acids Res* 2017, **45**:D491-D498.
- 1022 76. Tatusov RL, Galperin MY, Natale DA, Koonin EV: **The COG database: a tool**
1023 **for genome-scale analysis of protein functions and evolution.** *Nucleic Acids*
1024 *Res* 2000, **28**:33-36.
- 1025 77. Huerta-Cepas J, Szklarczyk D, Forslund K, Cook H, Heller D, Walter MC, Rattei
1026 T, Mende DR, Sunagawa S, Kuhn M, et al: **eggNOG 4.5: a hierarchical**

- 1027 **orthology framework with improved functional annotations for eukaryotic,**
1028 **prokaryotic and viral sequences. *Nucleic Acids Res* 2016, **44**:D286-293.**
- 1029 78. Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M: **KEGG for integration**
1030 **and interpretation of large-scale molecular data sets. *Nucleic acids research***
1031 **2012, **40**:D109-114.**
- 1032 79. Marchler-Bauer A, Derbyshire MK, Gonzales NR, Lu S, Chitsaz F, Geer LY, Geer
1033 RC, He J, Gwadz M, Hurwitz DI, et al: **CDD: NCBI's conserved domain**
1034 **database. *Nucleic Acids Res* 2015, **43**:D222-226.**
- 1035 80. Finn RD, Coghill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC,
1036 Punta M, Qureshi M, Sangrador-Vegas A, et al: **The Pfam protein families**
1037 **database: towards a more sustainable future. *Nucleic Acids Res* 2016,**
1038 ****44**:D279-285.**
- 1039 81. Bru C, Courcelle E, Carrere S, Beausse Y, Dalmar S, Kahn D: **The ProDom**
1040 **database of protein domain families: more emphasis on 3D. *Nucleic Acids***
1041 ***Res* 2005, **33**:D212-215.**
- 1042 82. Sigrist CJ, Cerutti L, Hulo N, Gattiker A, Falquet L, Pagni M, Bairoch A, Bucher
1043 P: **PROSITE: a documented database using patterns and profiles as motif**
1044 **descriptors. *Brief Bioinform* 2002, **3**:265-274.**
- 1045 83. Haft DH, Selengut JD, White O: **The TIGRFAMs database of protein families.**
1046 ***Nucleic Acids Res* 2003, **31**:371-373.**
- 1047 84. Hunter S, Apweiler R, Attwood TK, Bairoch A, Bateman A, Binns D, Bork P, Das
1048 U, Daugherty L, Duquenne L, et al: **InterPro: the integrative protein signature**
1049 **database. *Nucleic acids research* 2009, **37**:D211-215.**
- 1050 85. Sillitoe I, Lewis TE, Cuff A, Das S, Ashford P, Dawson NL, Furnham N,
1051 Laskowski RA, Lee D, Lees JG, et al: **CATH: comprehensive structural and**
1052 **functional annotations for genome sequences. *Nucleic Acids Res* 2015,**
1053 ****43**:D376-381.**
- 1054 86. Lam SD, Dawson NL, Das S, Sillitoe I, Ashford P, Lee D, Lehtinen S, Orengo
1055 CA, Lees JG: **Gene3D: expanding the utility of domain assignments. *Nucleic***
1056 ***Acids Res* 2016, **44**:D404-409.**
- 1057 87. Letunic I, Doerks T, Bork P: **SMART: recent updates, new developments and**
1058 **status in 2015. *Nucleic Acids Res* 2015, **43**:D257-260.**
- 1059 88. Potenza E, Di Domenico T, Walsh I, Tosatto SC: **MobiDB 2.0: an improved**
1060 **database of intrinsically disordered and mobile proteins. *Nucleic Acids Res***
1061 **2015, **43**:D315-320.**
- 1062 89. Pedruzzi I, Rivoire C, Auchincloss AH, Coudert E, Keller G, de Castro E, Baratin
1063 D, Cuhe BA, Bougueleret L, Poux S, et al: **HAMAP in 2015: updates to the**
1064 **protein family classification and annotation system. *Nucleic Acids Res* 2015,**
1065 ****43**:D1064-1070.**
- 1066 90. Nikolskaya AN, Arighi CN, Huang H, Barker WC, Wu CH: **PIRSF family**
1067 **classification system for protein functional and evolutionary analysis. *Evol***
1068 ***Bioinform Online* 2007, **2**:197-209.**
- 1069 91. Attwood TK, Beck ME, Flower DR, Scordis P, Selley JN: **The PRINTS protein**
1070 **fingerprint database in its fifth year. *Nucleic Acids Research* 1998, **26**:304-**
1071 **308.**

- 1072 92. Gene Ontology C: **The Gene Ontology project in 2008.** *Nucleic Acids Res*
1073 2008, **36**:D440-444.
- 1074 93. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment**
1075 **search tool.** *Journal of Molecular Biology* 1990, **215**:403-410.
1076 [<http://blast.wustl.edu>]
- 1077 95. Notredame C, Higgins DG, Heringa J: **T-Coffee: A novel method for fast and**
1078 **accurate multiple sequence alignment.** *J Mol Biol* 2000, **302**:205-217.
- 1079 96. Langmead B, Salzberg SL: **Fast gapped-read alignment with Bowtie 2.** *Nat*
1080 *Methods* 2012, **9**:357-359.
- 1081 97. Coordinators NR: **Database Resources of the National Center for**
1082 **Biotechnology Information.** *Nucleic Acids Res* 2017, **45**:D12-D17.
- 1083 98. Jari Oksanen FGB, Michael Friendly, Roeland Kindt, Pierre Legendre, Dan
1084 McGlenn, Peter R. Minchin, R. B. O'Hara, Gavin L. Simpson, Peter Solymos, M.
1085 Henry H. Stevens, Eduard Szoecs and Helene Wagner **vegan: Community**
1086 **Ecology Package. R package.** version 2.4-1. edition; 2016.
- 1087 99. **ggsignif** [<https://github.com/const-ae/ggsignif>]
- 1088 100. Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'Hara RB, Simpson
1089 GL, Solymos P, Stevens MHH, Wagner H: **vegan: Community Ecology**
1090 **Package.** *R package version 20-2* 2011.
- 1091 101. Cohen A: **On the graphical display of the significant components in a two-**
1092 **way contingency table.** *Communications in Statistics—Theory and Methods*
1093 1980, **A9**:1025–1041.
- 1094 102. Friendly M: **Graphical methods for categorical data.** *SAS User Group*
1095 *International Conference Proceedings* 1992, **17**:190–200.
- 1096 103. Zeileis A, Meyer D, Hornik K: **Residual-based Shadings for Visualizing**
1097 **(Conditional) Independence.** *Journal of Computational and Graphical Statistics*
1098 2007, **16**:507-525.
1099
- 1100

1101 **Declarations**

1102 Electronic supplementary material

1103

1104 **Additional file 1: Figure S1.** Boxplot of the proportion of sequencing reads after
1105 removing human contaminants from the samples between different Community State
1106 Types (CSTs). CSTs were defined as previously according to the composition and
1107 structure of the microbial community [29]. Plotted are interquartile ranges (IQRs, boxes),
1108 medians (line in box), and mean (red diamond). Significance value was calculated using
1109 Wilcoxon rank sum test using *ggsignif* R package [99]. Star sign (*) denotes the level of
1110 significance.

1111

1112 **Additional file 1: Figure S2.** Heatmap of relative abundance of the 50 most abundant
1113 phylotypes in the vaginal metagenomes used in this study. Ward linkage clustering is
1114 used to clusters samples based on their Jensen-Shannon distance calculated in the
1115 *vegan* package in R [100] according to the previous naming convention [29]. The
1116 sidebars indicate CSTs and gene richness category, respectively. Gene richness
1117 categories include high gene count (HGC) and low gene count (LGC), defined using the
1118 threshold of 10,000 genes per sample.

1119

1120 **Additional file 1: Figure S3.** Vaginal community accumulation curves and diversity
1121 estimate. (A) Accumulative diversity estimates with respect to sample size, for rarefied
1122 and extrapolated estimates using all samples; (B) accumulative diversity estimates with
1123 respect to sample size, for rarefied and extrapolated estimates using samples of
1124 different CSTs; (C) diversity estimate with respect to sample coverage, for rarefied and
1125 extrapolated estimates using all samples; (D) diversity estimate with respect to sample
1126 coverage, for rarefied and extrapolated estimate using samples of different CSTs.
1127 Community diversity estimates were computed using R package *iNEXT* [71] and *vegan*
1128 [72]. Sampling curve was either rarefied to smaller sample sizes or extrapolated to a
1129 larger sample size for species diversity estimate.

1130

1131 **Additional file 1: Figure S4.** Pie chart taxonomic distribution of reads that failed to map
1132 on VIRGO for vaginal metagenomes of African women from Gosmann *et al.* [30] in **A**
1133 and of Chinese women from [31] in **B**. The unmapped reads were compared to
1134 GenBank nt database [97] using BLASTN.

1135
1136 **Additional file 1: Figure S5.** Pipeline for data processing and integration for the
1137 construction of the human vaginal integrated non-redundant gene catalogue (VIRGO)
1138 and vaginal orthologous protein family groups (VOG). Metagenomes from 264 vaginal
1139 metagenomes and 416 genomes of urogenital isolates were processed, that including
1140 212 *in-house* sequenced vaginal metagenomes. The procedures include pre-processing
1141 to remove human contaminates, quality assessment, metagenome assembly, gene
1142 calling, functional and taxonomic annotation, gene clustering based on nucleotide
1143 sequencing similarity to form VIRGO, and Jaccard index coefficient clustering of
1144 amino acid sequences to form VOG.

1145
1146 **Additional file 1: Figure S6.** Proportion of the assembly length assigned taxonomically
1147 from the samples (**A**) among different community state types (CSTs) and (**B**) between
1148 different gene richness category. CSTs were defined as previously according to the
1149 composition and structure of the microbial community [29]. Gene richness category
1150 includes high gene count (HGC) and low gene count (LGC), defined using the threshold
1151 of 10,000 genes per sample.

1152
1153 **Additional file 1: Figure S7.** Top 20 species with the most abundant gene content in
1154 VIRGO. The ratio of the gene content of a species over the entire community to the
1155 base 2. Plotted are interquartile ranges (IQRs, boxes), medians (line in box), and mean
1156 (red diamond).

1157
1158 **Additional file 1: Figure S8.** Boxplot of the alignment scores of Jaccard orthologous
1159 clusters (JOCs) with multiple members. The alignment program T-Coffee [95] was used
1160 to access the alignment quality using alignment score.

1161

1162 **Additional file 1: Figure S9.** Phylogeny that is demonstrative use of VOG to
1163 characterize the *G. vaginalis* cholesterol-dependent cytolysin (CDC) protein family. It
1164 shows the phylogeny of CDC-containing protein and alignment of domain 4 of the CDCs
1165 that is generally well conserved but contains a single divergent site, highlighted in
1166 yellow [38].

1167

1168 **Additional file 1: Figure S10.** Association plot of functional distribution of different
1169 gene count categories in vaginal microbiome. Functional category was defined using
1170 EggNOG (v4.5) [77] functional category. A Cohen-Friendly association plot [101, 102]
1171 was produced in statistical package *vcd* in R [103] to indicate deviations to indicate
1172 deviations from independence of CSTs and functional distribution. Mosaics display was
1173 shown, where the cells are shaded in proportion to standardized residuals, where the
1174 positive value (blue) is the observed frequency is substantially greater than would be
1175 found under independence, and the negative value (red) indicates cells which occur
1176 less often than under independence.

1177

1178 **Additional file 1: Figure S12.** Functional category of *L. iners* in different gene richness
1179 categories. Functional category was defined using EggNOG (v4.5) [77] functional
1180 category.

1181

1182 **Additional file 1: Figure S13.** Gene richness category and taxonomic distribution of
1183 tryptophan production-related genes in VIRGO. **(A)** Pie chart of the percentage of
1184 tryptophan production-related genes in different gene richness categories of HGC or
1185 LGC. **(B)** The top 10 most affiliated taxonomic groups of the tryptophan production-
1186 related genes.

1187

1188 **Additional file 1: Figure S14.** Heatmap includes gene prevalence profiling of available
1189 genomes of vaginal isolates and VIRGO-characterized metagenomes for **(A)** *L.*
1190 *crispatus*, **(B)** *L. iners*, **(C)** *L. jensenii*, **(D)** *L. gasseri*, and **(E)** *G. vaginalis*, **(F)** *A. vaginae*
1191 and **(G)** *P. timonensis*. Hierarchical clustering of the profiles was performed using ward

1192 linkage based on Jaccard similarity coefficient. CSTs were defined as previously
1193 according to the composition and structure of the microbial community [29].

1194

1195 **Additional file 2: Table S1.** Statistics of the sequence reads, including 211 *in-house*
1196 sequenced metagenomes, 53 metagenomes from HMP DACC database, 277 genomes
1197 isolated from vagina, reproductive or urinary system deposited in GenBank and 139
1198 urogenital bacteria genomes from HMP DACC database used to compile database. The
1199 assembly statistics includes assembled base pairs, total number of contigs, N50, mean
1200 and median length, and other statistics.

1201

1202 **Additional file 2: Table S2.** OTUs table for all metagenomes included in VIRGO.
1203 Taxonomic profiling was conducted in MetaPhlAn version 2 [69]. Community state types
1204 were defined as previously according to the composition and structure of the microbial
1205 community [29]. 312 bacterial species that were present in at $\geq 0.01\%$ relative
1206 abundance are shown.

1207

1208 **Additional file 2: Table S3.** Statistics of the complete and subsets of the sequence
1209 contigs included into VIRGO, including reference data sets: i) complete VIRGO
1210 database, ii) 212 *in-house* sequenced vaginal metagenomes, iii) 53 HMP DACC vaginal
1211 metagenomes [32], iv) all HMP urogenital reference genomes, v) 277 genomes of
1212 bacteria isolated from vagina, reproductive or urinary system deposited in GenBank,
1213 and vi) 139 genomes of urogenital bacteria from HMP DACC database [15].

1214

1215 **Additional file 2: Table S4.** Table showing counts of the non-redundant genes in
1216 VIRGO by taxonomic groups in both species and genera.

1217

1218 **Additional file 2: Table S5.** The vaginal fungal database that includes 5 vaginal yeast
1219 species in 40 genomes and the abundance of detected fungal and phage in
1220 metagenome samples.

1221

1222 **Additional file 2: Table S6.** Annotation and alignment of a Jaccard orthologous
1223 clusters (JOCs) involved in vaginolysin. This JOCs was in one protein family in VOG
1224 that contains multiple genes, annotation information is in **A**. Multiple sequence
1225 alignment of this protein family was performed in T-Coffee [95], was used to access the
1226 alignment quality (**B**).

1227

1228 **Additional file 2: Table S7.** Examples of cell surface-associated proteins of *L. iners*.
1229 Two Jaccard orthologous clusters (JOCs) involved in this function were retrieved from
1230 VIRGO. (**A**) the JOC was recognized to have LPXTG motif; (**B**) the JOC that harbor
1231 motif YSIRK.

1232

1233 **Additional file 2: Table S8.** The statistics of number of non-redundant genes identified
1234 in a metagenome and the depth sequencing for samples in different CSTs.

1235

1236 **Additional file 2: Table S9.** Examples of tryptophan production-related gene cataloging
1237 using VIRGO. It includes three essentials genes Tryptophanase (*TnaA*), Tryptophan
1238 synthase beta chain (*TrpB*), and Tryptophanyl-tRNA synthetase (*TrpS*) in gene name,
1239 gene richness, gene annotation, to demonstrate the profiling of a specific function of
1240 interest and its taxonomic distribution.

1241

1242 **Additional file 2: Table S10.** Summary of the 7 vaginal bacterial species with gene
1243 content characterized using VIRGO to determine the diversity of individual populations.
1244 It includes four *Lactobacillus* species (*L. crispatus*, *L. iners*, *L. jensenii*, and *L. gasseri*),
1245 as well as three additional species common to the vagina (*G. vaginalis*, *A. vaginae* and
1246 *P. timonensis*). Reads mapping was performed using 1,507 *in-house* and publicly
1247 available vaginal metagenomes to VIRGO. Metagenomes that contained at least 80% of
1248 their average genome's number of coding genes were included. Abbr: Av: *A. vaginae*;
1249 Gv: *G. vaginalis*; Pt: *P. timonensis*; Lc: *L. crispatus*; Li: *L. iners*; Lj: *L. jensenii*; Lg: *L.*
1250 *gasseri*.

1251

1252 **Additional file 2: Table S11.** List of accession numbers for genomes of the four
1253 *Lactobacillus* species including *L. crispatus*, *L. iners*, *L. jensenii*, and *L. gasseri* and
1254 three species including *G. vaginalis*, *A. vaginae* and *P. timonensis* used in intraspecies
1255 analyses.

1256

1257 **Additional file 2: Table S12.** Taxonomic distribution of pullulanase domain-containing
1258 proteins included in VIRGO.

1259

1260 Availability of data and material

1261 All database data and code were made freely assessable on [https://github.com/Ravel-](https://github.com/Ravel-Laboratory/VIRGO)
1262 [Laboratory/VIRGO](https://github.com/Ravel-Laboratory/VIRGO). It includes Jaccard index clustering code, VIRGO non-redundant
1263 nucleotide gene database, VOG amino acid protein family database, curated taxonomy
1264 and functions information, and tutorials. Metagenomes used in the analyses are
1265 deposited at EBA ### (The list of accession numbers for the 1,507 vaginal
1266 metagenomes used in intraspecies analyses will be available upon acceptance of the
1267 manuscript).

1268

1269 Competing interests

1270 The authors declare no competing interests.

1271 Authors' contributions

1272 B.M., J.R. designed the research. B.M., M.F., J.H., and J.R. performed the research.
1273 B.M., M.H. generated the data. B.M., M.F., J.H., and J.C. analyzed the data. B.M., M.F.,
1274 R.B., and J.R. interpreted the data and wrote the paper.

1275 Acknowledgements

1276 The authors thank Drs. Douglas Kwon and Matthew Hayward for their helpful
1277 assistance in analyzing African women metagenomes. The authors thank Dr. Nan Qin
1278 and Qian Xu for their helpful assistance in analyzing Chinese women metagenomes.
1279 Research reported in this publication was supported by the National Institutes of Allergy
1280 and Infectious Diseases and Nursing Research of the National Institutes of Health

1281 under award numbers U19AI084044, R01NR015495 and R01AI116799, and the Bill &
1282 Melinda Gates Foundation award OPP1189217.

1283 **Figures**

1284

1285 **Figure 1.** Percent of vaginal metagenome reads that can be mapped to contigs from the
1286 following reference data sets: i) complete VIRGO database, ii) 211 *in-house* sequenced
1287 vaginal metagenomes, iii) 53 HMP DACC vaginal metagenomes [32], iv) all HMP
1288 urogenital reference genomes, v) 277 genomes of bacteria isolated from vagina,
1289 reproductive or urinary system deposited in GenBank, and vi) 139 genomes of
1290 urogenital bacteria from HMP DACC database [15]. Values plotted are the average,
1291 error bars represent the standard error of the mean.

1292

1293 **Figure 2.** Pipeline for data processing and integration for the construction of the human
1294 vaginal integrated non-redundant gene catalog (VIRGO) and vaginal orthologous
1295 groups (VOG) for protein families. Metagenomes from 264 vaginal metagenomes and
1296 416 genomes of urogenital isolates were processed, that including 211 *in-house*
1297 sequenced vaginal metagenomes. The procedures include preprocessing to remove
1298 human contaminants, quality assessment, metagenome assembly, gene calling,
1299 functional and taxonomic annotation, gene clustering based on nucleotide sequencing
1300 similarity to form VIRGO, and jaccard index coefficient clustering of amino acid
1301 sequences to form VOG. A more detailed illustration is in **Additional file 1: Figure S5**
1302 and description is in Material and Method.

1303

1304 **Figure 3.** Taxonomic and functional composition of vaginal microbiome in VIRGO. **(A)**
1305 Top 20 species with the most abundant gene content in VIRGO. The logarithm of the
1306 ratio of the gene content of a species over the entire community to the base 2. Plotted
1307 are interquartile ranges (IQRs, boxes), medians (line in box), and mean (red diamond).
1308 **(B)** Species-specific metagenome accumulation curves for the number of non-
1309 redundant genes. **(C)** Functional distribution of non-redundant genes in VIRGO.
1310 Functional categories were defined using EggNOG (v4.5) [77]. **(D)** Prevalence of
1311 BVAB1 in metagenomes using a minimum number of genes threshold of 50% of the
1312 estimated BVAB1 genome size. A gene was present if ≥ 3 reads mapped to it. **(E)**
1313 Relationship between the depth of sequencing and the number of bacterial non-

1314 redundant genes identified using VIRGO. Each point is a separate metagenome and is
1315 color-coded according to community state type.

1316

1317 **Figure 4. (A)** Boxplot of the number non-redundant genes in samples of different
1318 Community State Types (CSTs). CSTs were defined as previously according to the
1319 composition and structure of the microbial community [29]. Table below boxplot
1320 contains percentage of samples in each of the CSTs stratified by high gene count
1321 (HGC) or low gene count (LGC). **(B)** Plot of the \log_2 transformed ratio of the gene of a
1322 species being in one gene count category over the other across the 264 vaginal
1323 metagenomes, only the species with more than 4 times more abundant in a category
1324 (either HGC or LGC) are shown. Plotted are interquartile ranges (IQRs, boxes),
1325 medians (line in box), and mean (red diamond).

1326

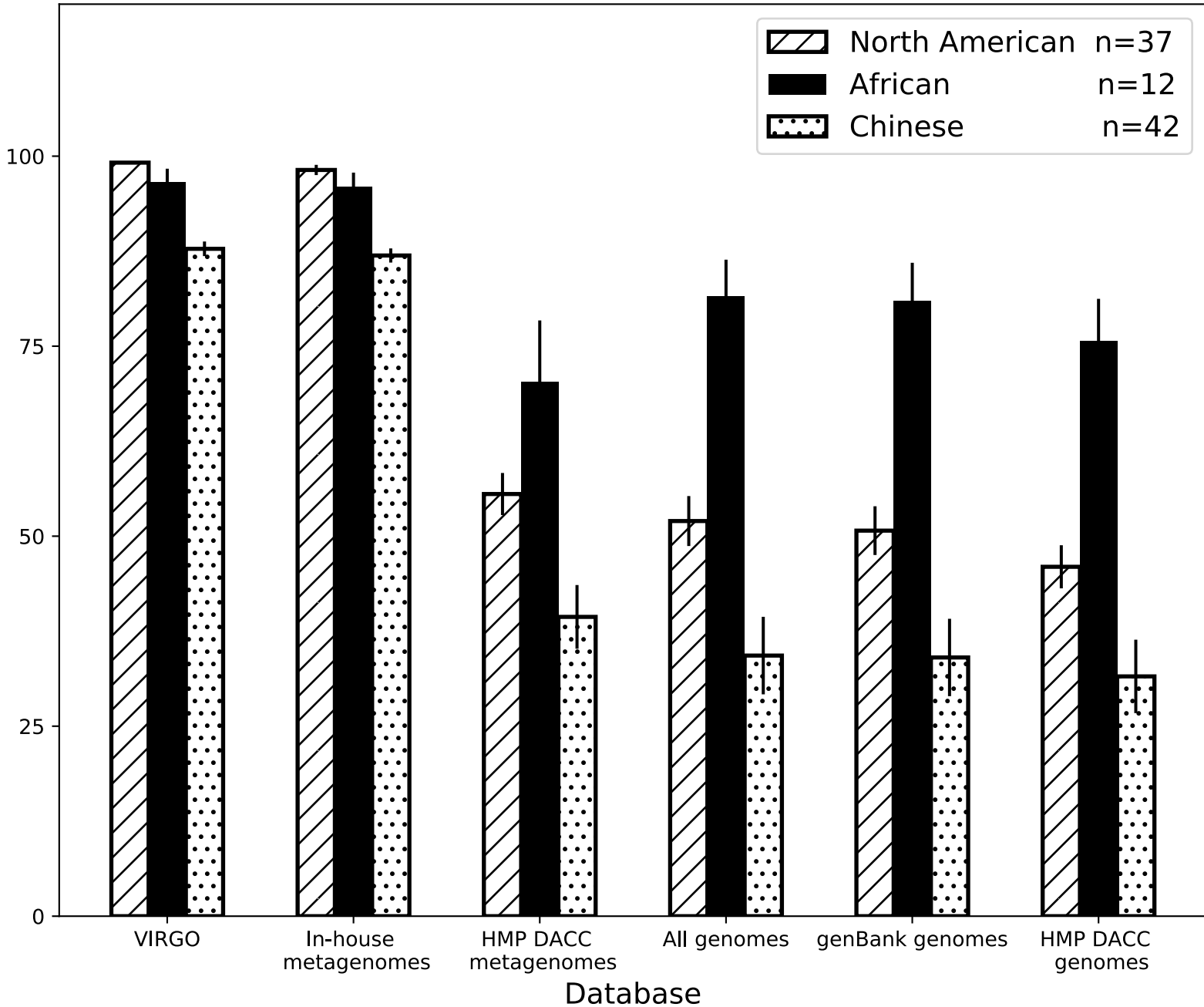
1327 **Figure 5.** Demonstration using VIRGO and VOG to study vaginal microbiome. **(A)** 4
1328 sampling points were selected based on a longitudinally profiled subject prior to (T1),
1329 during (T2 and T3), and after (T4) an episode of bacterial vaginosis using 16S rRNA
1330 profiling. **(B)** Functional profiling of the metagenome (MG) and metatranscriptome (MT)
1331 of each of the 4 sampling points. Functional categories were annotated using EggNOG
1332 (v4.5) [77]. **(C)** Functional profiles stratified by species using the taxonomic profiling
1333 provided by VIRGO. **(D)** Demonstrative use of VOG to characterize the *G. vaginalis*
1334 cholesterol-dependent cytolysin (CDC) protein family. It shows the phylogeny of CDC-
1335 containing protein and alignment of domain 4 of the CDCs that is generally well
1336 conserved but contains a single divergent site, highlighted in yellow [38].

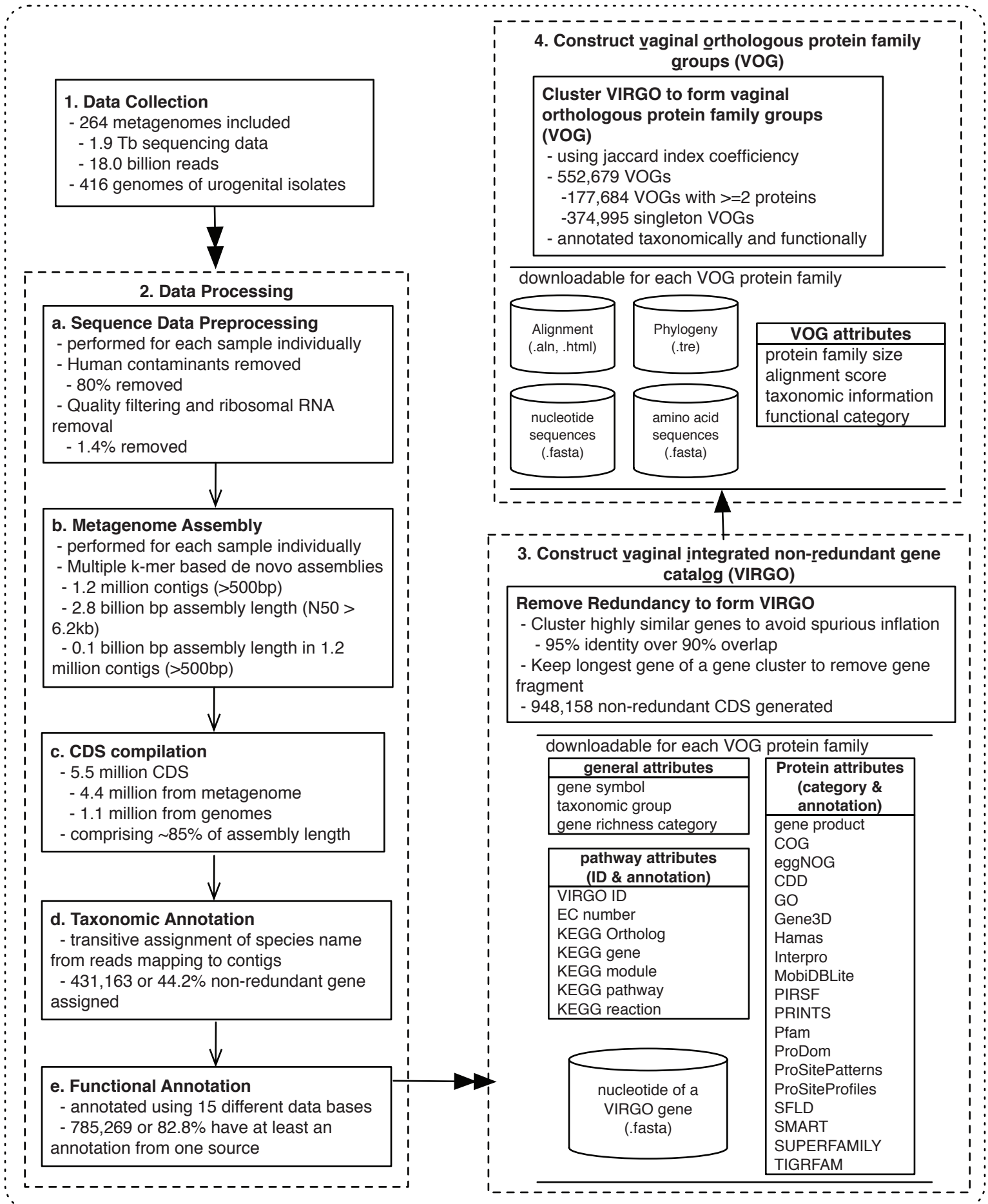
1337

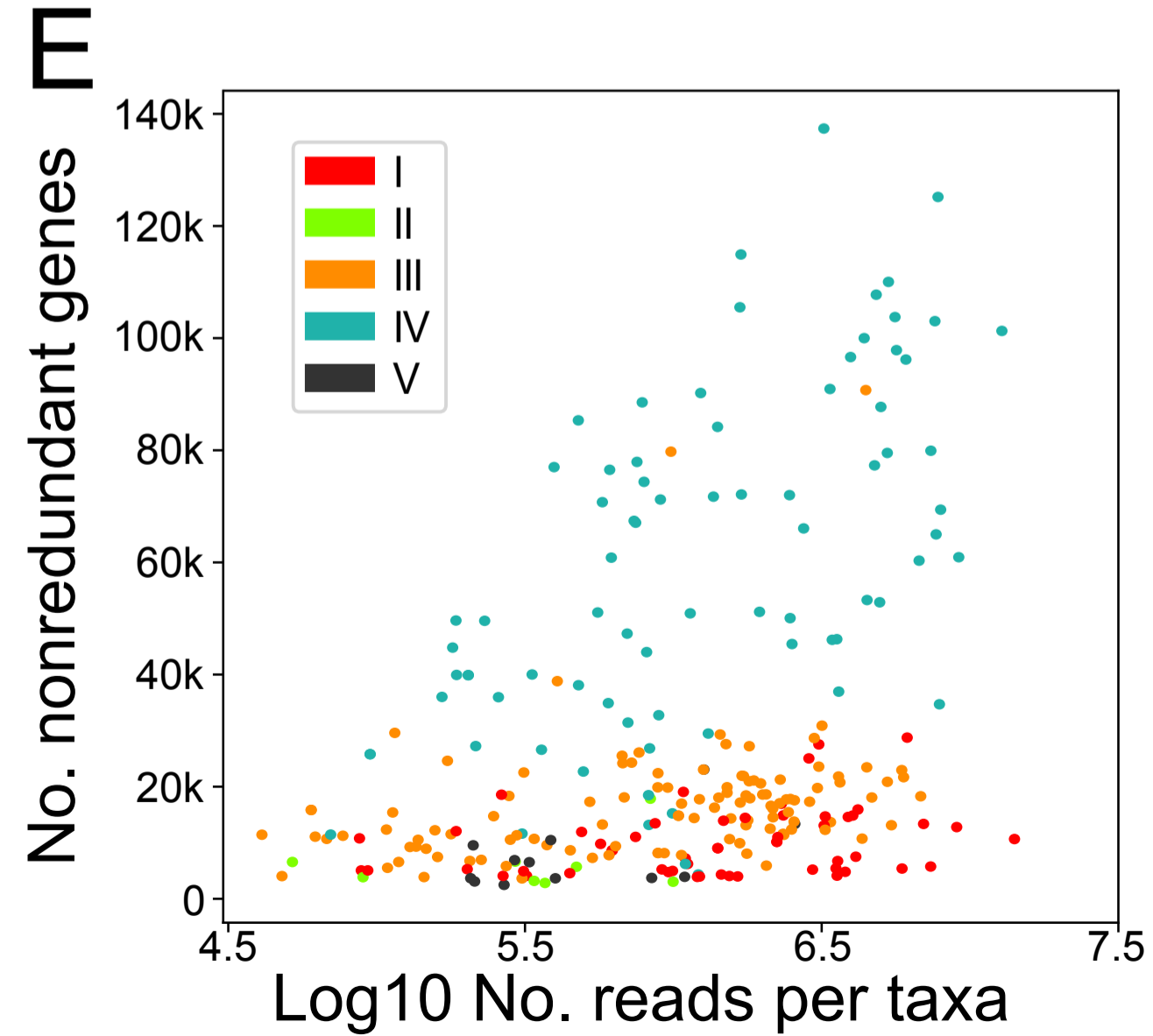
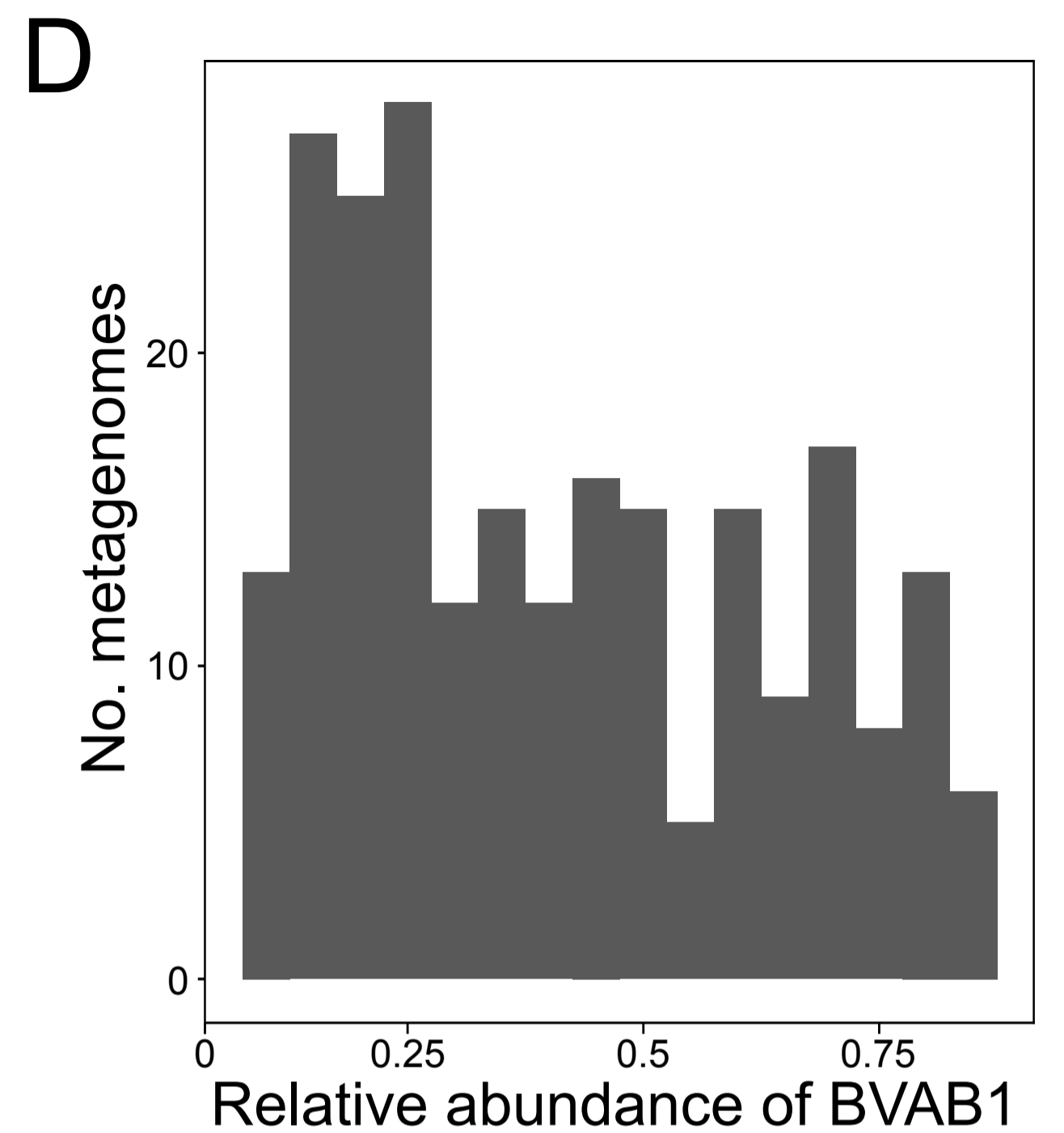
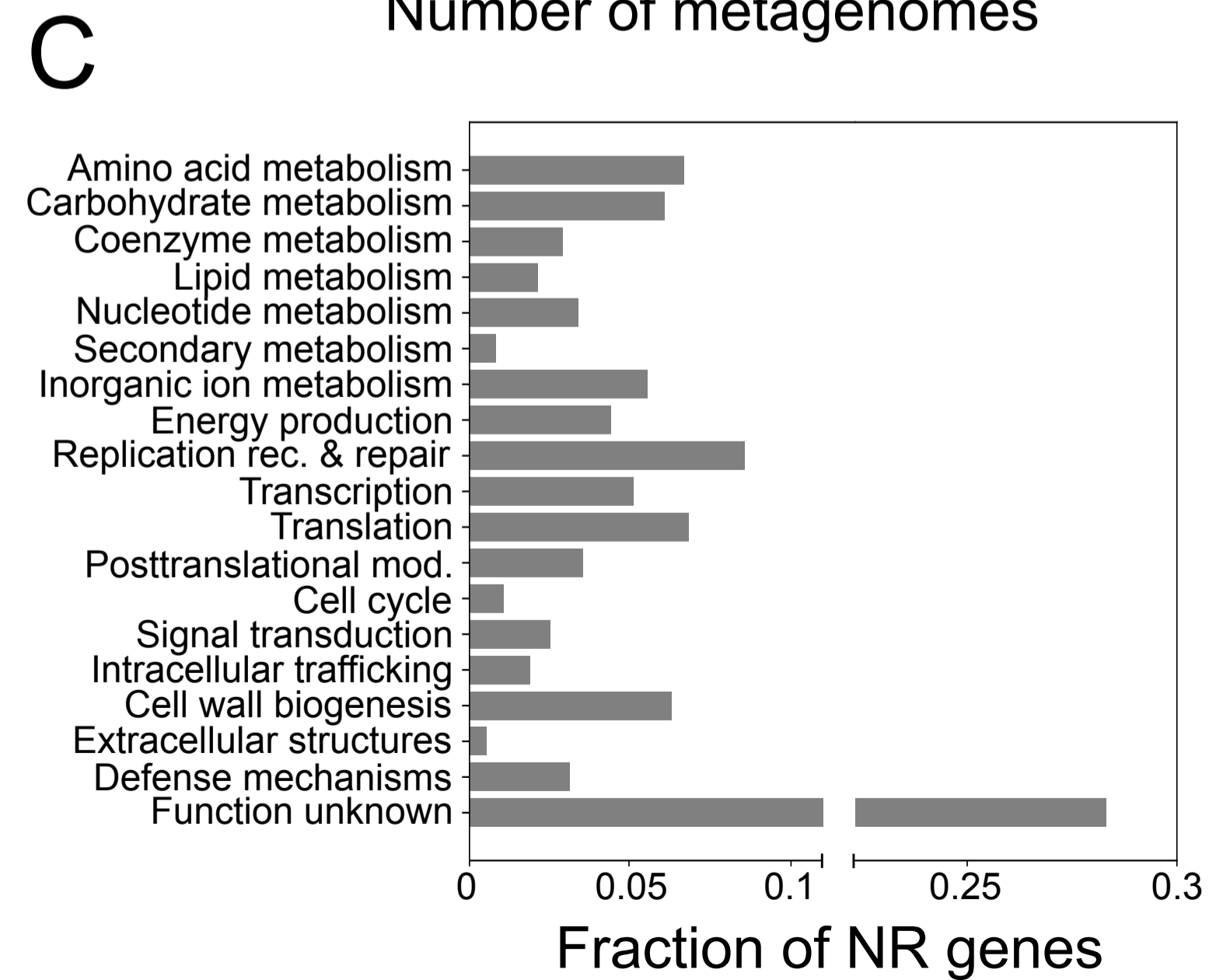
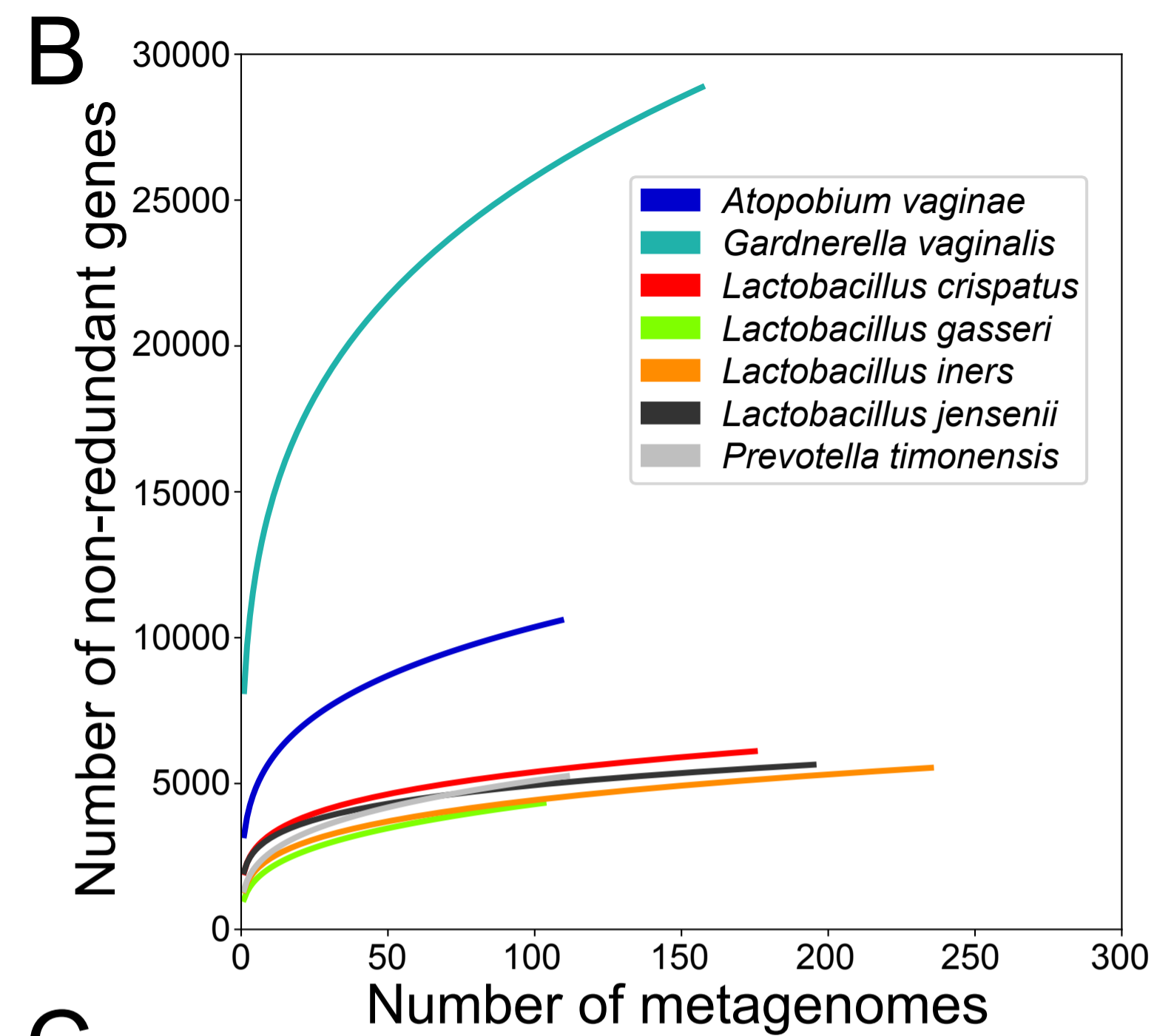
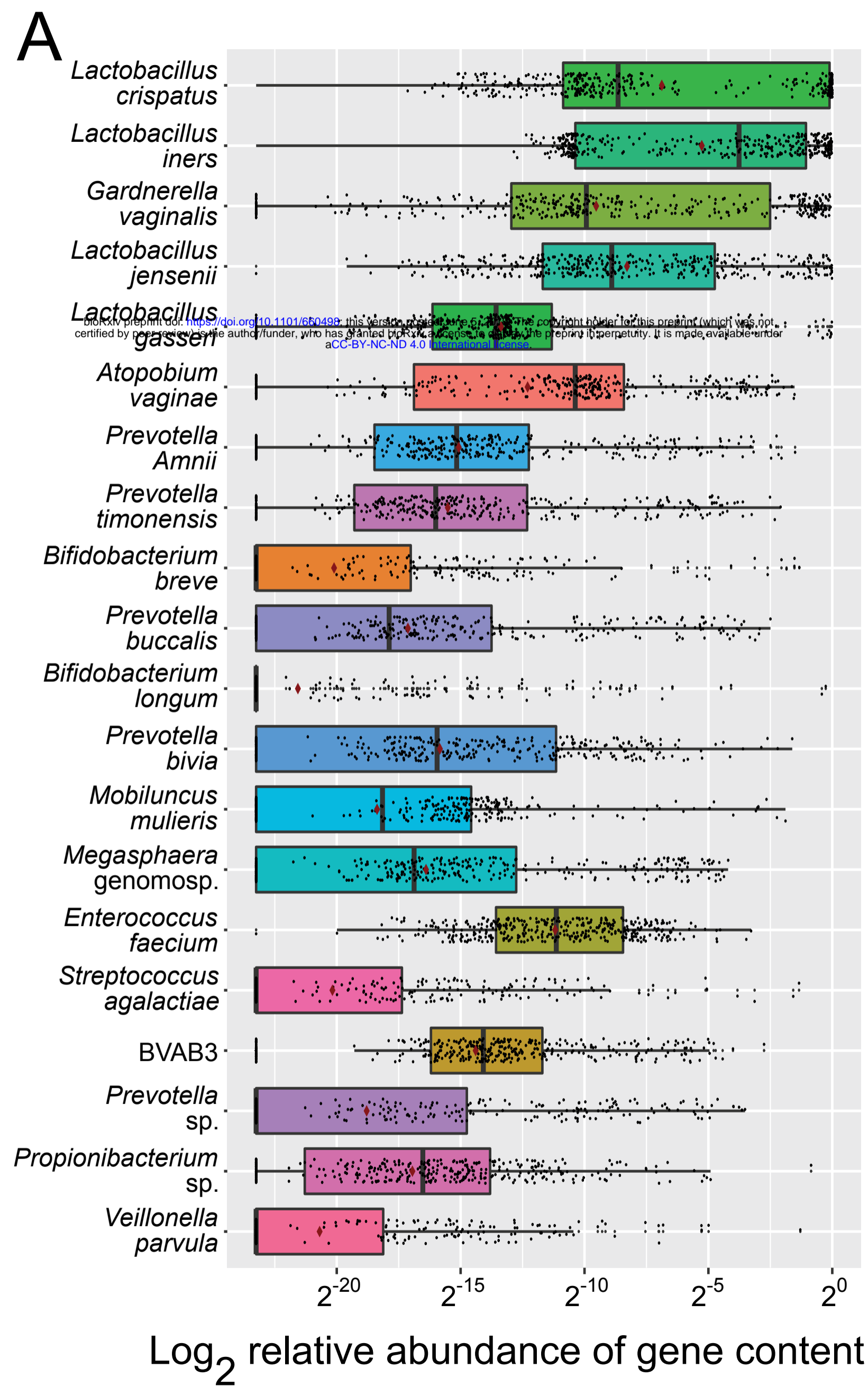
1338 **Figure 6.** Intraspecies diversity revealed using VIRGO of seven vaginal species
1339 including *L. crispatus*, *L. iners*, *L. jensenii*, *L. gasseri*, and *G. vaginalis*, *A. vaginae* and
1340 *P. timonensis*. **(A)** Summary of the number (N) of isolate genomes and metagenome
1341 (MG) samples with more than 80% of their average genome's number of coding genes
1342 for a species, based on a dataset of 1,507 *in-house* vaginal metagenomes
1343 characterized using VIRGO. **(B)** Boxplot of number non-redundant genes in isolate
1344 genomes versus vaginal metagenomes. **(C)** Heatmap of presence/absence of *L.*

1345 *crispatus* non-redundant gene profiles for 56 available isolate genomes and 413
1346 VIRGO-characterized metagenomes that contained either high (red) or low (blue)
1347 relative abundance of the species. Hierarchical clustering of the profiles was performed
1348 using ward linkage based on their Jaccard similarity coefficient. * number of isolate
1349 genomes and metagenome samples. † MG: Metagenomes * $p < 0.05$, *** $p < 0.001$ after
1350 correction for multiple comparisons.

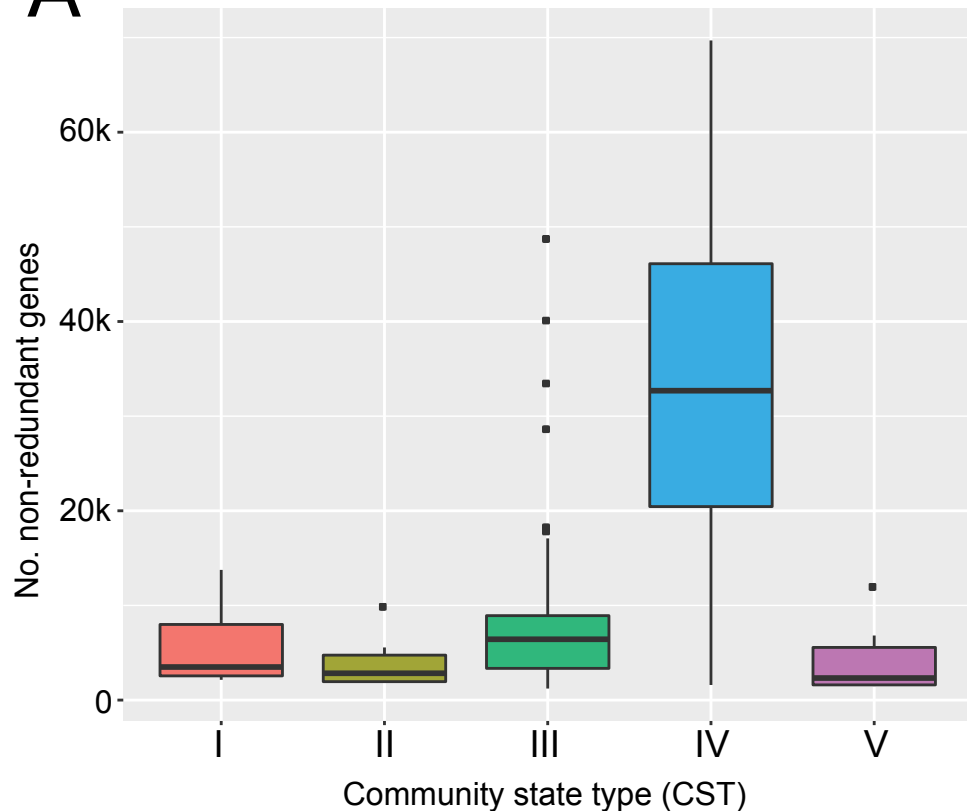
Percent of reads mapped to database







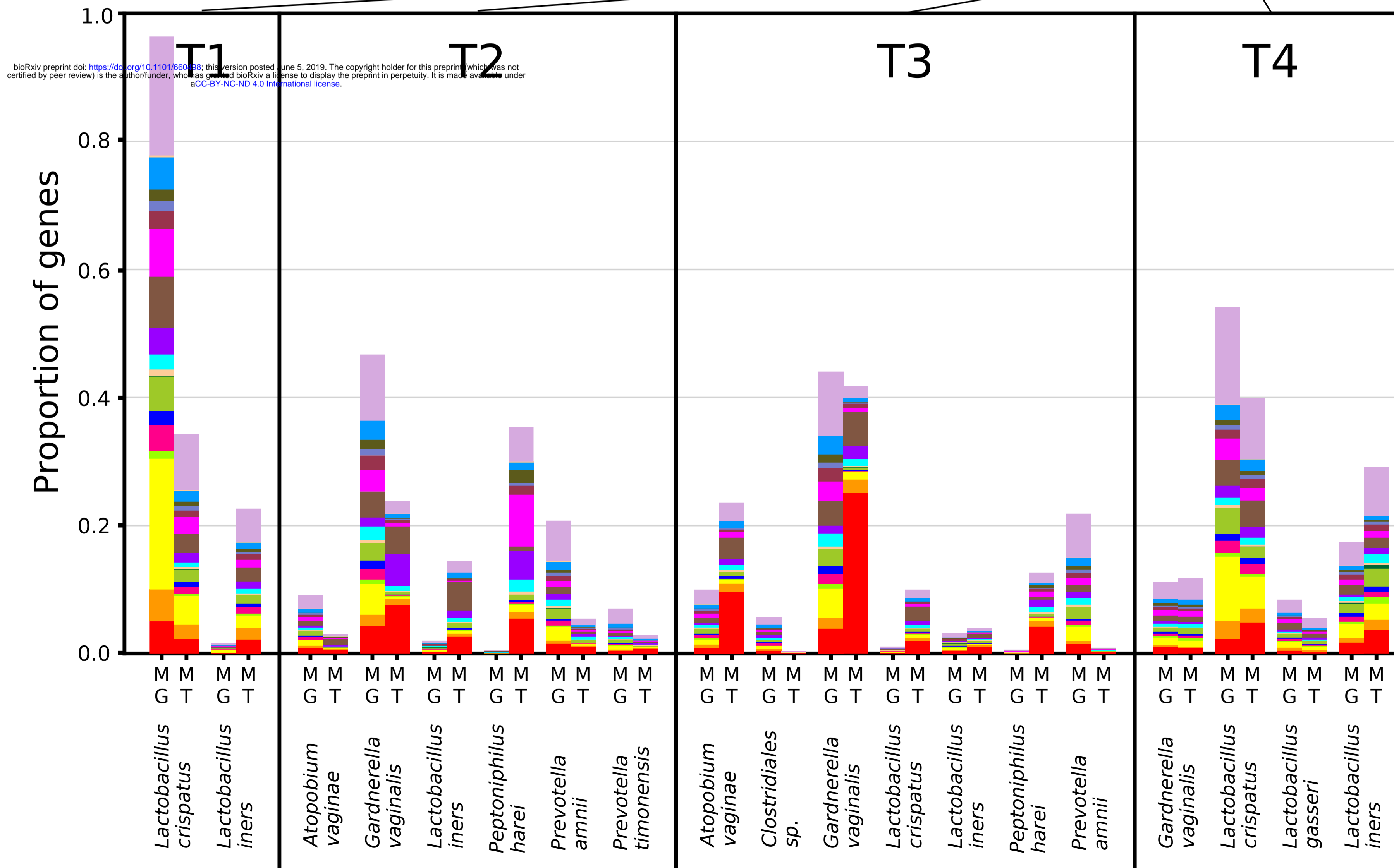
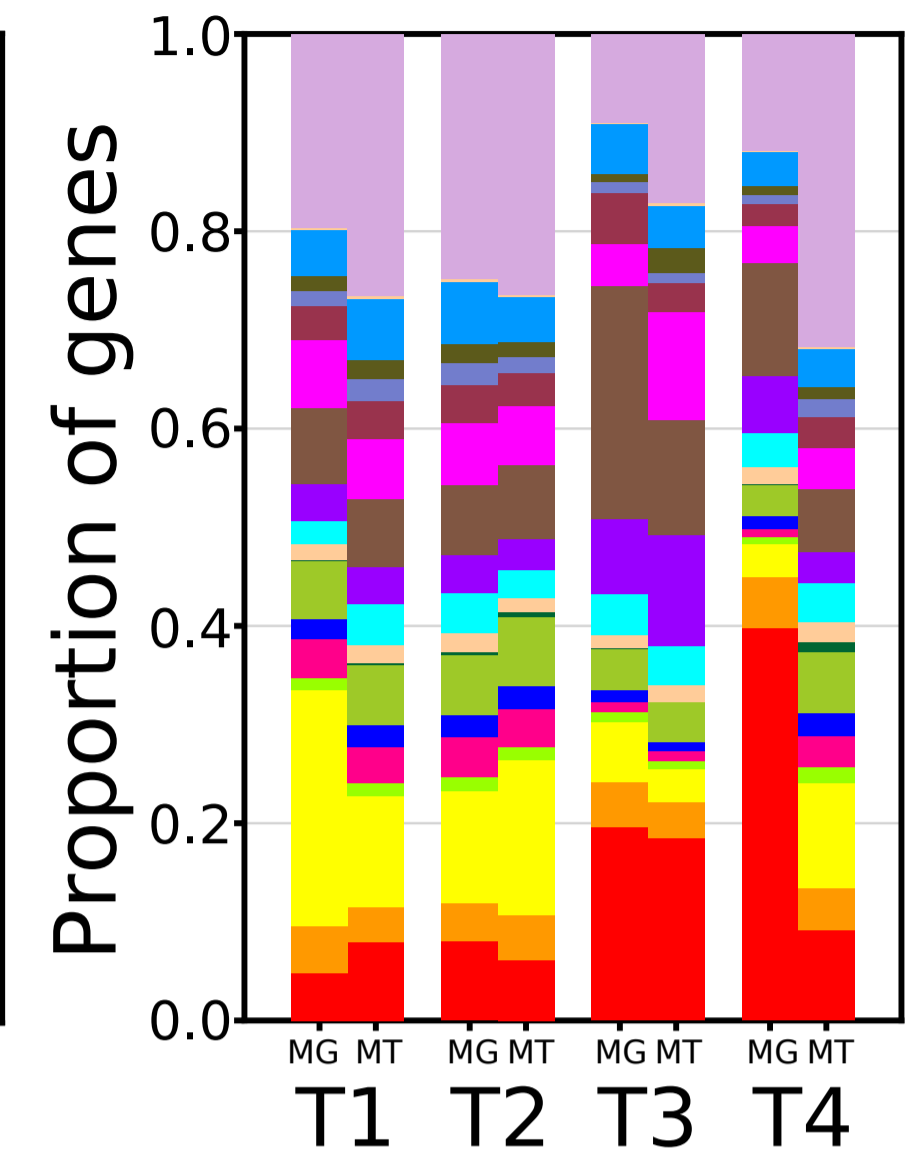
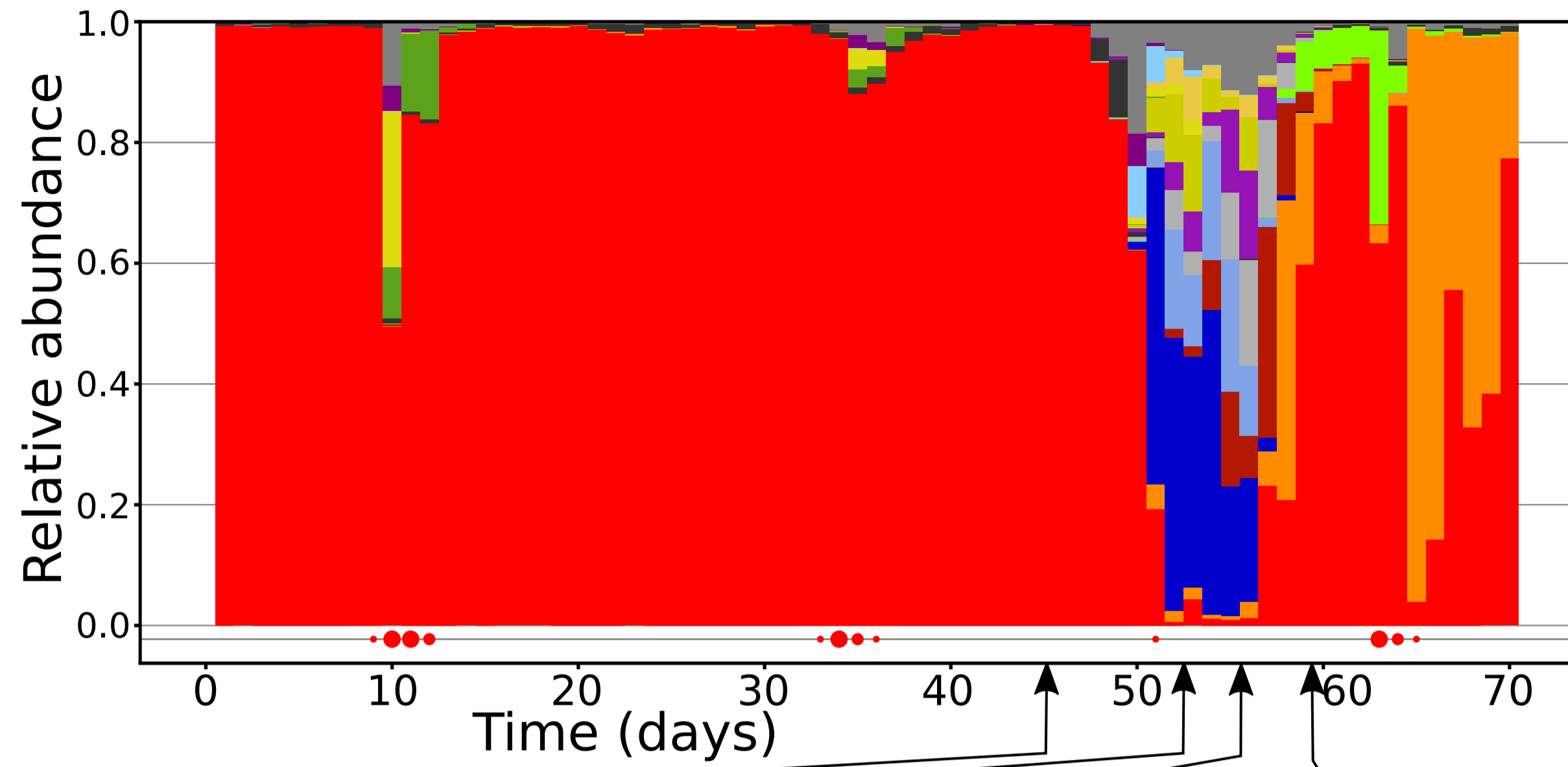
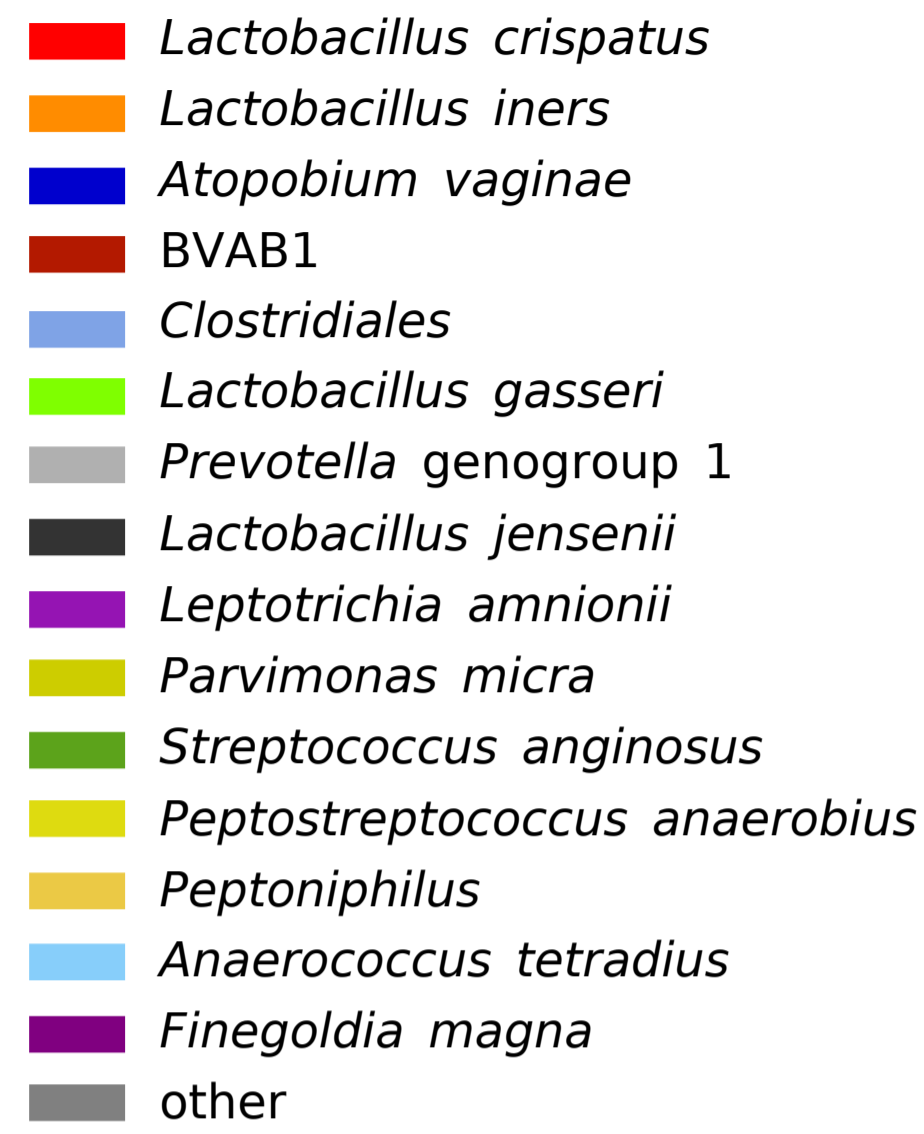
A



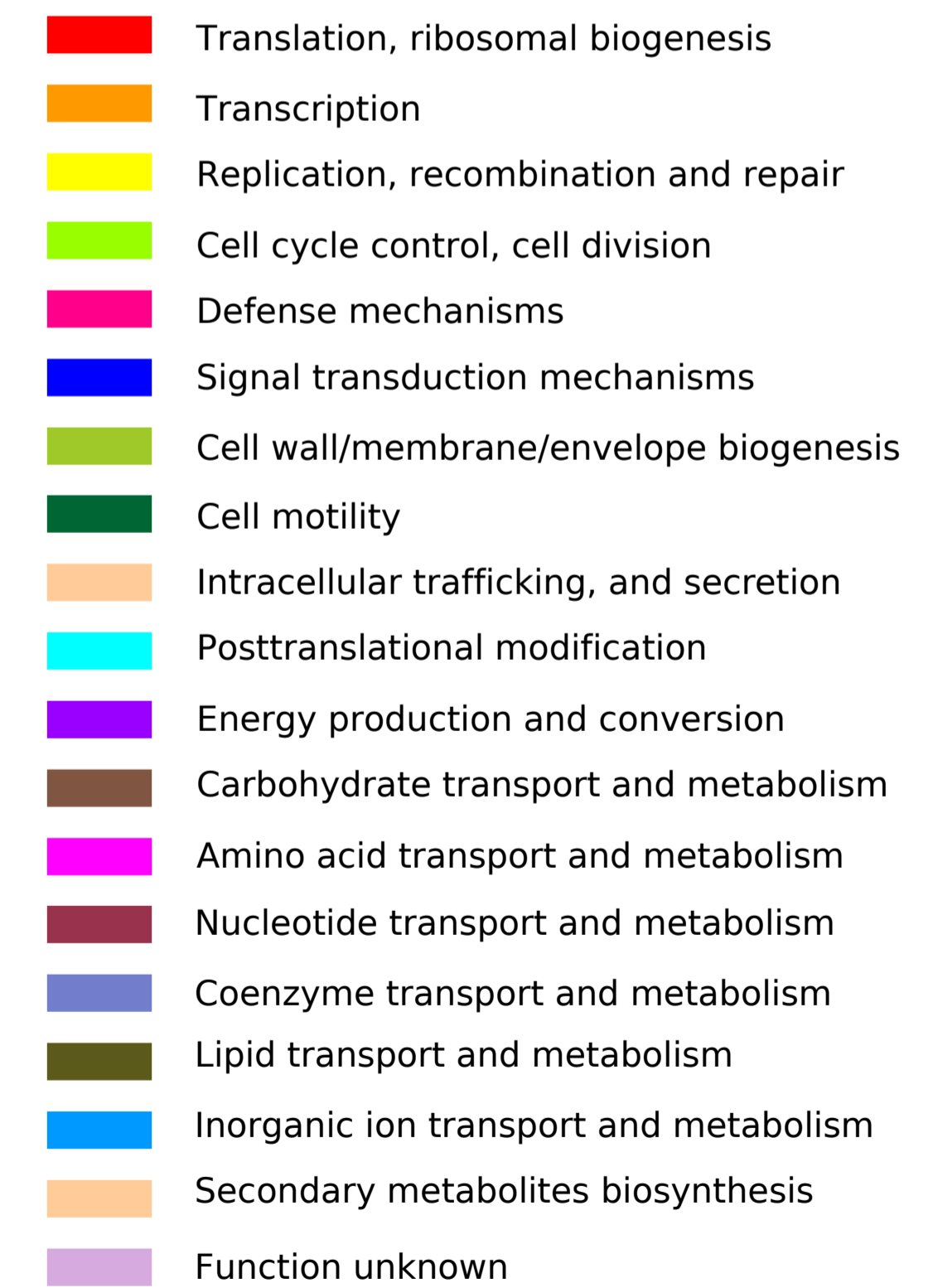
| CST | I | II | III | IV | V |
|-----|-------|------|-------|-------|-------|
| LGC | 88.5% | 100% | 78.3% | 11.7% | 91.7% |
| HGC | 11.5% | 0% | 21.7% | 88.3% | 8.3% |

B



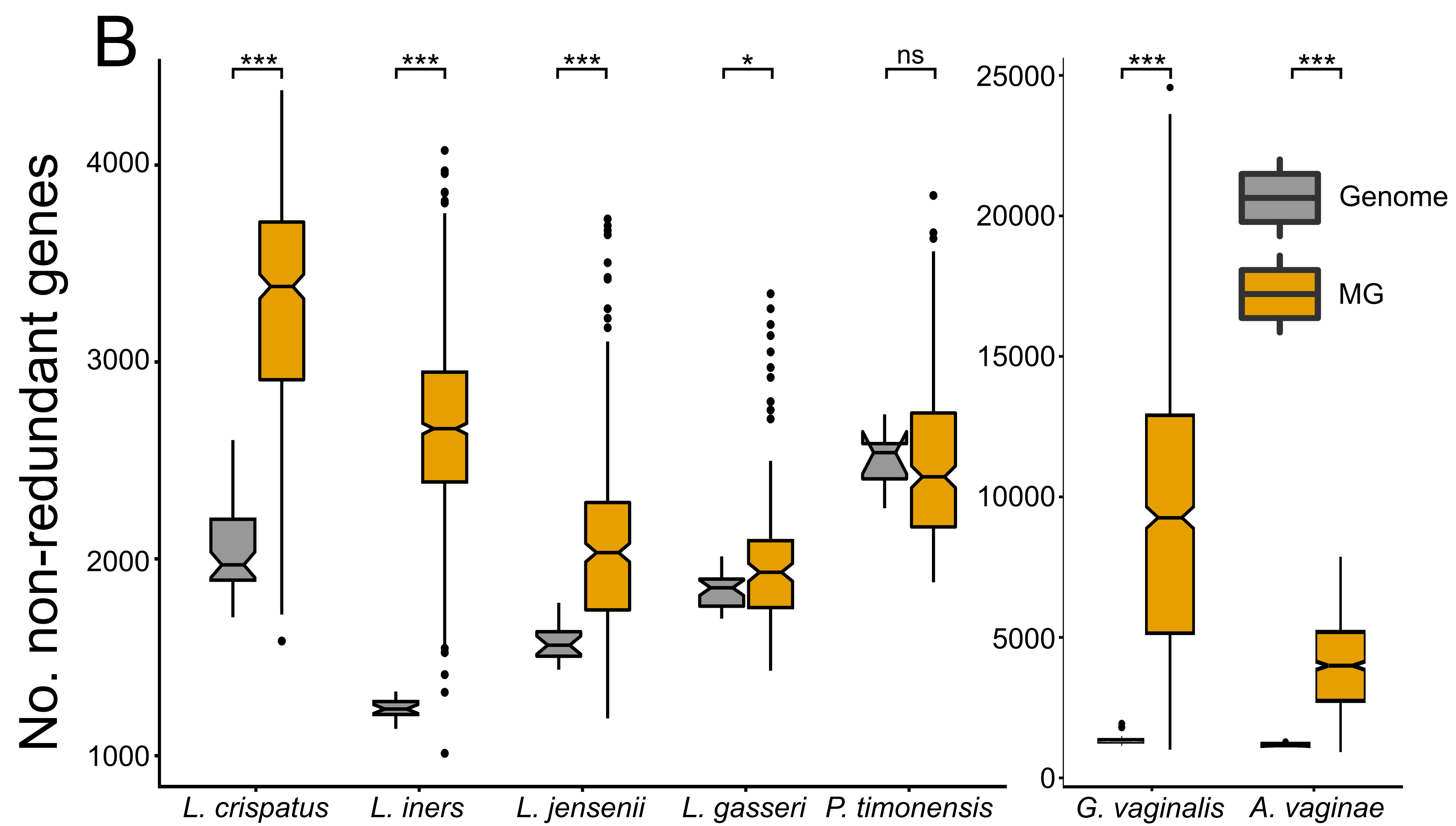


COG category



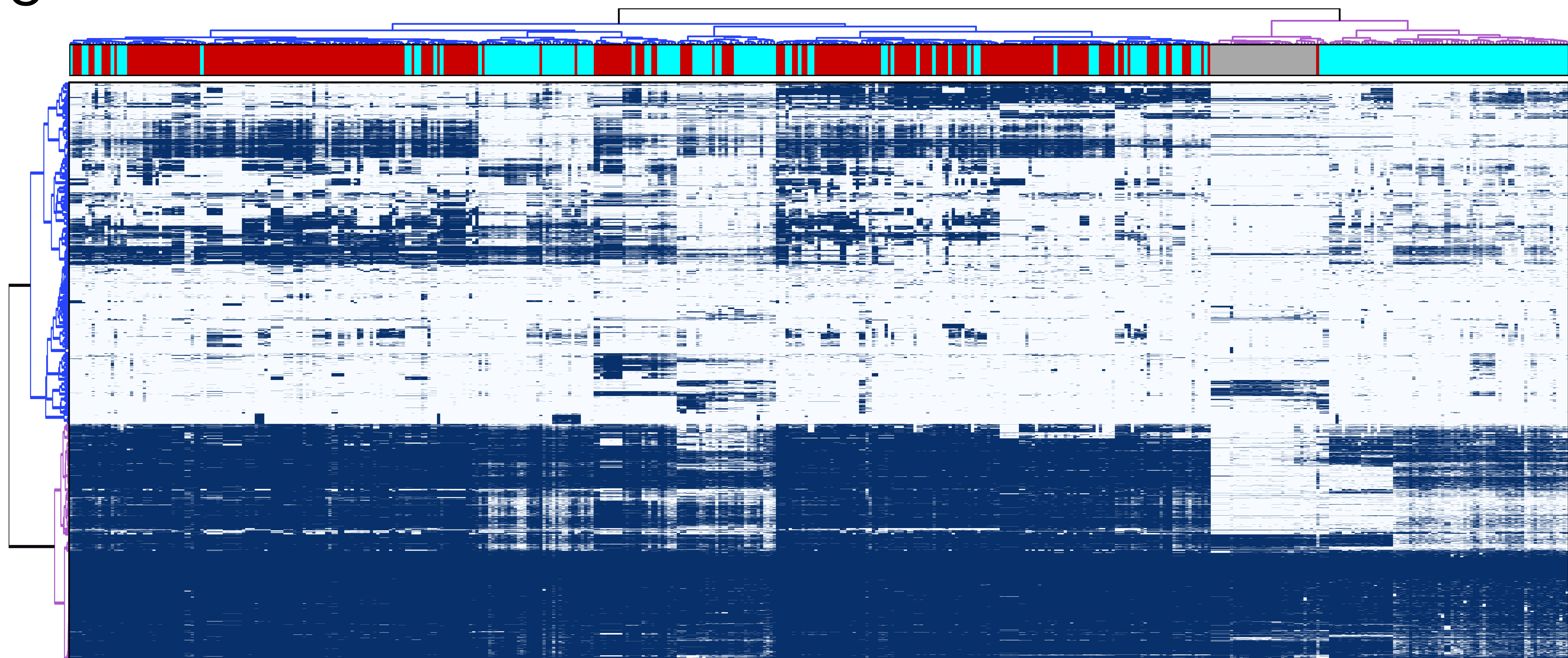
A

| Species | Sequence type | N | Ave. No. genes | SD |
|----------------------|---------------|------|----------------|------|
| <i>L. crispatus</i> | Genome | 56 | 2064 | 225 |
| | MG | 413 | 3262 | 568 |
| <i>L. iners</i> | Genome | 21 | 1236 | 45 |
| | MG | 1028 | 2659 | 458 |
| <i>L. jensenii</i> | Genome | 18 | 1566 | 89 |
| | MG | 334 | 2041 | 471 |
| <i>L. gasseri</i> | Genome | 31 | 1830 | 89 |
| | MG | 139 | 1989 | 371 |
| <i>P. timonensis</i> | Genome | 7 | 2498 | 158 |
| | MG | 273 | 2469 | 401 |
| <i>G. vaginalis</i> | Genome | 90 | 1329 | 122 |
| | MG | 1007 | 9301 | 5079 |
| <i>A. vaginae</i> | Genome | 6 | 1188 | 41 |
| | MG | 724 | 4007 | 1658 |



C Metagenomes containing *L. crispatus* & genomes from isolates

L. crispatus non-redundant genes



D Metagenomes containing *L. gasseri* & genomes from isolates

L. gasseri non-redundant genes

