

1 **Resolution and co-occurrence patterns of *Gardnerella leopoldii*, *Gardnerella***
2 ***swidsinskii*, *Gardnerella piovii* and *Gardnerella vaginalis* within the vaginal**
3 **microbiome**

4

5 Janet E. Hill^{1,*}, Arianne Y.K. Albert² and the VOGUE Research Group

6

7 ¹Department of Veterinary Microbiology, University of Saskatchewan, Saskatoon, SK S7N
8 5B4, Canada, Janet.Hill@usask.ca

9 ²Women's Health Research Institute, BC Women's Hospital and Health Centre,
10 Vancouver, BC V6H 3N1, Canada, Arianne.Albert@cw.bc.ca

11

12 *To whom correspondence should be addressed

13

14 JEH: Janet.Hill@usask.ca

15 AYKA: Arianne.Albert@cw.bc.ca

16

17 JEH ORCID: [0000-0002-2187-6277](https://orcid.org/0000-0002-2187-6277)

18 AYKA ORCID: [0000-0002-3050-0888](https://orcid.org/0000-0002-3050-0888)

19

20 **Abstract**

21 **Background**

22 *Gardnerella vaginalis* is a hallmark of vaginal dysbiosis, but is found in the microbiomes
23 of women with and without vaginal symptoms. *G. vaginalis* encompasses diverse taxa
24 differing in attributes that are potentially important for virulence, and there is evidence that
25 ‘clades’ or ‘subgroups’ within the species are differentially associated with clinical
26 outcomes. The *G. vaginalis* species description was recently emended, and three new
27 species within the genus were defined (*leopoldii*, *swidsinskii*, *piotii*). 16S rRNA sequences
28 for the four *Gardnerella* species are all >98.5% identical and no signature sequences
29 differentiate them.

30 **Results**

31 We demonstrated that *Gardnerella* species can be resolved using partial chaperonin-60
32 (cpn60) sequences, with pairwise percent identities of 87.1-97.8% among the type strains.
33 Pairwise co-occurrence patterns of *Gardnerella* spp. in the vaginal microbiomes of 413
34 reproductive aged Canadian women were investigated, and several significant co-
35 occurrences of species were identified. Abundance of *G. vaginalis*, and *swidsinskii* was
36 associated with vaginal symptoms of abnormal odour and discharge.

37 **Conclusions**

38 cpn60 barcode sequencing can provide a rapid assessment of the relative abundance of
39 *Gardnerella* spp. in microbiome samples, providing a powerful method of elucidating
40 associations between these diverse organisms and clinical outcomes. Researchers should
41 consider using cpn60 in place of 16S RNA for better resolution of these important
42 organisms.

43 **Background**

44 Since its original isolation from human vaginal samples in 1953 [1], the species that
45 eventually became known as *Gardnerella vaginalis* has been strongly associated with
46 vaginal dysbiosis and negative reproductive outcomes [2]. Evaluation of the abundance of
47 *Gardnerella* morphotypes in Gram stained vaginal smears is a key factor in calculation of
48 the Nugent score [3] and thus in the microbiological definition of bacterial vaginosis (BV).
49 Understanding of its role in the vaginal microbiome has been complicated due to
50 phenotypic diversity within the taxon and its perplexing presence, sometimes in high
51 numbers, in the vaginal microbiomes of women without any signs or symptoms of
52 dysbiosis [2].

53 Over the past decades, several classification schemes have been developed in an
54 attempt to delineate subgroups within *Gardnerella vaginalis*. These have included both
55 “biotyping” schemes based on a set of biochemical test results [4, 5], and molecular
56 methods based on amplification and restriction digestion of 16S rRNA gene sequences [6].
57 In 2005, four clusters of *G. vaginalis*-like sequences were observed in vaginal microbiome
58 profiles based on sequencing of amplified cpn60 barcode sequences [7]. More recently,
59 four “subgroups” or “clades” of *G. vaginalis* were defined based either on partial sequences
60 of the cpn60 barcode sequence [8-10] or the concatenated sequences of 473 genes common
61 to a set of 17 *G. vaginalis* isolates [11]. When we compared these latter two approaches
62 directly, they were consistent with each other, with cpn60 defined subgroups A-D
63 corresponding to clades 4, 2, 1 and 3, respectively [8]. Based on the results of this
64 comparison, we suggested that pairwise average nucleotide identity (ANI) values between
65 whole genome sequences of representative isolates were consistent with their definition as

66 separate species [12] but that additional phenotypic characteristics that differentiate the
67 subgroups should be identified [8].

68 In 2019, Vaneeccoutte et al. [13] formally proposed the emendment of the species
69 *G. vaginalis*, and defined three new species: *G. piotti*, *G. swidsinskii* and *G. leopoldii*, based
70 on comparison of whole genome sequences, biochemical properties and matrix-assisted
71 laser desorption ionization time-of-flight (MALDI-TOF) mass spectrometry analysis. The
72 species could not be resolved using 16S rRNA gene sequences [13]. In addition to these
73 four species, the authors also defined nine additional “genome species” based on whole
74 genome sequence comparisons. These additional genome species were not named or
75 formally described, presumably due to a lack of sufficient numbers of isolates to make a
76 strong case for their designation.

77 Given the apparent significance of *Gardnerella* spp. in the vaginal microbiome,
78 resolution of these species in metagenomic samples and association of their presence and
79 abundance with clinical outcomes is critical. We have already demonstrated that
80 amplification and sequencing of the cpn60 barcode can be used to resolve cpn60-defined
81 subgroups A-D in vaginal samples, and demonstrated that their differential abundance can
82 be used to reveal previously unreported community state types [14]. The objectives of the
83 current study were to determine if cpn60 barcode sequences could differentiate the newly
84 defined species and genome species of *Gardnerella*, to investigate their distribution in a
85 collection of 417 previously sequenced vaginal microbiome profiles, and to identify
86 associations of *Gardnerella* spp. with vaginal symptoms. Our results confirm the resolving
87 power of the cpn60 barcode sequence and reveal significant co-occurrences of *Gardnerella*

88 spp. in the vaginal microbiome that have implications for diagnostics for women's health,
89 and for our understanding of vaginal microbial ecology.

90

91 **Methods**

92 ***Gardnerella cpn60 sequence analysis***

93 cpn60 universal target sequences from 52 *Gardnerella* species representing the four
94 named species (*G. vaginalis*, *G. piotti*, *G. swidsinskii* and *G. leopoldii*) and nine additional
95 genome species described by Vaneechoutte et al. [13] were retrieved from cpnDB
96 (www.cpnadb.ca, [15, 16]) and aligned using CLUSTALw. Inter-species nucleotide
97 sequence similarities were calculated using *dnadist* within PHYLIP [17]. A bootstrapped
98 phylogenetic tree was calculated using *seqboot*, *dnadist* (maximum likelihood option), and
99 *neighbor*. The type strain of *Alloscardovia omnicolens* (DSM 21503) was included as an
100 outgroup. The consensus tree was computed with *consense* and branch lengths applied with
101 *fitch*. The tree was visualized using FigTree (v1.4.2).

102 ***Identification of Gardnerella spp. in vaginal microbiomes of Canadian women***

103 cpn60 barcode sequence data, Nugent scores and self-reported symptom data from
104 previously conducted studies by the VOGUE Research Group of the vaginal microbiome
105 composition of non-pregnant, reproductive aged Canadian women recruited from clinics
106 in greater Vancouver, Canada area were used in the current sub-analysis. These included
107 healthy women (n=310), women living with HIV (n=54) and women who had at least four
108 self-identified episodes of vulvovaginitis in the past 12 months (n=53). DNA extraction,
109 cpn60 barcode PCR, library preparation and sequencing of amplicons are described in [14].
110 Amplification primer sequences were removed using cutadapt, followed by quality

111 trimming with trimmomatic (quality cut-off 30, minimum length 150). Quality trimmed
112 reads were loaded into QIIME2 [18] for sequence variant calling and read frequency
113 calculation with DADA2 [19] and a truncation length of 250. For taxonomic identification,
114 variant sequences were compared to the cpnDB_nr reference database (version 20190305,
115 downloaded from www.cpnadb.ca) using watered-BLAST [20]. The reference database
116 includes representative sequences of each of the four named *Gardnerella* species and the
117 nine additional genome species defined by Vaneechoutte et al. [13].

118 ***Co-occurrence and cluster analysis***

119 For the following analyses we removed four samples with total sequence reads
120 <500, leaving 413 samples. We determined pairwise co-occurrences using the presence or
121 absence of each *Gardnerella* species or genome species in each sample. Presence was
122 indicated if the raw number of sequence reads in the sample was ≥ 10 , otherwise the species
123 was marked as absent. We calculated a Jaccard index of similarity (J) [21] for each pairwise
124 combination and compared these to a probabilistic null model of species co-occurrence that
125 takes into account observed frequencies to determine significance [22, 23](Supplemental
126 File 1). P-values were Benjamini-Hochberg (BH) corrected to a false-discovery rate = 0.05
127 [24]. Jaccard dissimilarities ($1-J$) were used to cluster the species using complete-linkage
128 hierarchical clustering implemented in the *vegan* package for R [25].

129 ***Comparisons of clinical data and relative abundance***

130 Relative abundances were calculated using the centre-log-ratio transformation as
131 implemented in the *ALDEx2* package [26] and *propr* package [29]. We compared clr
132 abundance among categories of clinical variables using Kruskal-Wallis tests. Significant
133 omnibus tests were followed up with Dunn's post-hoc tests with BH adjusted p-values. For

134 this analysis, the sample size was reduced to 395 for which we had concurrent Nugent
135 scoring data, and 393 for which we had self-report symptom data.

136

137 **Results and Discussion**

138 **Resolution of *Gardnerella* spp. based on cpn60 universal target sequences**

139 The four named *Gardnerella* spp. were resolved in the cpn60 phylogenetic tree
140 based on an alignment of the 552 bp “universal target” sequence barcode (Figure 1), with
141 good bootstrap support for nodes separating the species. *G. vaginalis* (genome sp. 1)
142 corresponds to the previously described subgroup C/clade 1. This observation is consistent
143 with previous demonstrations that cpn60 barcode sequences are generally excellent
144 predictors of whole genome sequence relationships among closely related bacteria [27, 28]

145 *G. swidsinskii* (genome sp. 6) and *G. leopoldii* (genome sp. 5) share a common
146 node in the tree, but are separated clearly with good bootstrap support. These two species
147 (represented by strains AMD and 5-1) were previously grouped together within cpn60
148 subgroup A/clade 4 based on the selection of isolates available for analysis at the time.
149 Vaneechoutte et al. noted in their description of the species definitions that *G. swidsinskii*
150 and *G. leopoldii* could be distinguished based on ANI, DNA-DNA hybridization (DDH)
151 and MALDI-TOF profiles [13].

152 *G. piovii* (genome sp. 4) corresponds to subgroup B/clade 2 and can be distinguished
153 from other species based on ANI, DDH, MALDI-TOF and a positive sialidase test. One *G.*
154 *piovii* isolate (JCP8151A) clustered with genome sp. 3 in the tree. In the original description
155 of the cpn60 subgroups, what was designated genome sp. 3 by Vaneechoutte et al. was
156 included in subgroup B and it is noteworthy that several isolates examined by Schellenberg

157 et al. have cpn60 sequences identical to strain 00703C2mash (genome sp. 3) were also
158 found to be sialidase positive [8]. The characterization of additional isolates of genome sp.
159 3 and *G. piotii* will be necessary to determine if these groups should be combined.

160 Subgroup D/clade 3 was the most diverse in previous descriptions and so it is not
161 surprising to find that isolates previously identified as Subgroup D (strains 101, 1500E,
162 6119V5 and 00703Dmash) are separated into three genome species: genome spp. 8, 9 and
163 10. Complete characterization and possible naming of these three genome species, along
164 with genome spp. 2, 7, and 11-13 will require analysis of additional isolates to establish
165 differentiation by whole genome sequence and additional phenotypic characteristics.

166 Pairwise nucleotide sequence identities for the type strains of *Gardnerella* spp. and
167 representatives of the other nine genome species were calculated from the aligned
168 sequences (Table S1). Identities among the four type strains were from 87.1% (*G. leopoldii*
169 vs. *G. piotii*) to 97.8% (*G. leopoldii* vs. *G. swidsinskii*). When representatives of the other
170 nine genome species were included, percent identities for the 552 bp cpn60 barcode
171 sequence ranged from 84.2% (genome sp. 7 vs. genome sp. 2 or *G. piotii*) to 99.4%
172 (genome sp. 9 vs. genome sp. 10). No isolates had identical cpn60 barcode sequences.
173 Inter-specific cpn60 barcode sequence identities are known to vary widely among bacteria
174 genera so this range was not unexpected [15]. To investigate the resolving power of the 5'
175 and 3' ends of the barcode sequence that would be determined using routine next-
176 generation sequencing protocols, we truncated the alignment to examine 250 bp of either
177 end of the barcode sequence. Average pairwise identities were 88.2% (range 83.2 - 99.6)
178 and 91.3% (range 84.4 - 99.2), respectively (Table S1). None of the species were identical.

179 These identities cover the same range as observed for the entire barcode sequence, as is
180 expected given the uniform distribution of sequence variation along its length [9].

181 **Classification of *Gardnerella* sequence variants**

182 One of the major advantages of use of the cpn60 barcode sequence for taxonomic
183 profiling of microbial communities is the ability to achieve species level classification of
184 sequence reads or assembled operational taxonomic unit (OTU) sequences routinely. It was
185 this resolution that led to the identification of previously undescribed community state
186 types in the human vaginal microbiome, based on the detection of subspecies level
187 sequences within *Gardnerella* [14]. Elucidation of the role of genomically and
188 phenotypically distinct *Gardnerella* lineages in the vaginal microbiome and determining
189 their association with clinical outcomes requires determining their distribution in clinical
190 cohorts. Accomplishing this on a large scale requires culture-independent techniques.
191 While whole-genome shotgun metagenomics might provide resolution of *Gardnerella*
192 spp., this approach requires orders of magnitude more sequencing effort and much more
193 complex bioinformatics than amplicon sequencing. Based on the successful resolution of
194 13 genome species of *Gardnerella* described above, we next sought to discover if they
195 could be reliably detected and quantified in cpn60 amplicon sequence-based microbiome
196 profiles.

197 cpn60 barcode sequence data was available from 417 previously characterized
198 vaginal samples from non-pregnant, reproductive aged Canadian women. For the purposes
199 of the current study, exact sequence variants were identified using DADA2 and a truncation
200 length of 250 bp and variants were compared to the cpnDB_nr reference database [15] to
201 identify the nearest database neighbour.

202 Most (301/413) of the samples for which at least 500 reads were available contained
203 some *Gardnerella*-like sequence variants and variants corresponding to all 13 *Gardnerella*
204 spp. and genome species were detected. The median sequence identity of variants to
205 reference sequences was 98.4%. Sample prevalence and proportional abundance ranged
206 widely among species (Table S2). For example, 68.4% (206/301) of *Gardnerella* positive
207 samples contained *G. vaginalis* and 49% (148/301) contained *G. swidsinskii*, but seven
208 genome species (2, 7-13) were detected in $\leq 10\%$ of samples. The prevalence and
209 abundance patterns are generally consistent with previous descriptions of vaginal
210 microbiomes based on cpn60 barcode sequencing [14, 29-32] or clade-specific quantitative
211 real-time PCR [33, 34]. There were 60 samples with at least 50% of their read counts
212 accounted for by *Gardnerella* spp., and 30 samples with at least 75% *Gardnerella* (Table
213 S2).

214 The number of *Gardnerella* spp. detected per sample ranged from one (109/301,
215 36.2%) to ten (3/301, 1%), although the majority (184/301, 61.1%) contained one or two
216 species (Figure 2). Overall, multiple *Gardnerella* spp. were detected in 63.8% of samples,
217 consistent with a previous report of multiple *Gardnerella* “clades” in 70% of samples from
218 women with BV [33]. The prevalence and proportional abundance of *Gardnerella* spp. in
219 samples with $>50\%$ *Gardnerella* ($n = 60$) are shown in Figure 3. In addition to the four
220 named species, genome sp. 3 was detected frequently and in relatively high proportional
221 abundance, in striking contrast to the rarely detected genome species 2, and 7-13.

222 The prevalence and abundance patterns we observed in these samples mirrors the
223 isolate and whole genome sequence collection used to provide evidence for the emendment
224 of *Gardnerella* and the designation of the new species [13]. Based on our experience, there

225 is no obvious bias in PCR amplification of *Gardnerella* lineages, and multiple
226 representatives of all previously defined cpn60 subgroups were readily amplified using
227 cpn60 “universal” PCR primers [8]. Furthermore, we have shown a strong correlation
228 between *Gardnerella* cpn60 sequence read counts in amplicon-based microbiome profiles
229 and abundances determined by *Gardnerella*-specific quantitative real-time PCR [29].
230 Thus, it seems that some *Gardnerella* species are actually less prevalent and do not achieve
231 proportional dominance in the populations of women we have examined to date.
232 Elucidating the ecological mechanisms responsible for this differential “success” of
233 *Gardnerella* spp. will require further focused study.

234 **Co-occurrence of *Gardnerella* spp.**

235 Given the frequency with which women are colonized by multiple species of
236 *Gardnerella*, we were interested to determine if there are any consistent patterns of co-
237 occurrence among species. Closely related species occupying similar environmental niches
238 might be expected to co-occur more frequently, and depending on resource levels, they
239 might also compete with each other. Raw correlations of read counts are not recommended
240 for assessing co-occurrence as they are biased in the context of the compositional nature of
241 amplicon sequence analysis [35-37]. Methods based on presence/absence, such as
242 Jaccard’s index can be informative in this context and perform better than raw correlations
243 [22]. To determine whether taxa co-occur more or less often than expected by chance, a
244 reasonable null model for Jaccard’s index is required. Traditional null models of co-
245 occurrence have used randomizations and simulations, but have been shown to be biased
246 under many circumstances [38]. Therefore, we used a probabilistic null model of co-
247 occurrence [23], which performs well for microbial sequencing data [22]. In addition to

248 co-occurrence analysis using presence/absence, we also investigated proportionality of
249 species [36] on the center-log ratio transformed read counts using the ‘propr’ package [39].
250 The results were very similar with *G. vaginalis* and *G. swidsinskii* showing clustering by
251 proportionality, as well as *G. piovii* and genome species 3. However, as there is currently
252 no agreed upon hypothesis testing method for proportionality, we used the
253 presence/absence data to determine significant co-occurrences.

254 Significant co-occurrences were observed for several pairs of species, but there
255 were also many cases where species co-occurred only randomly (Figure 4A, Table S3).
256 Among the most frequently detected species (*G. vaginalis*, *G. swidsinskii*, *G. leopoldii*, *G.*
257 *piovii* and genome sp. 3), the smallest pairwise Jaccard distances (i.e. the most samples in
258 common) were observed for *G. vaginalis* and *G. swidsinskii*, and *G. piovii* and genome sp.
259 3. (Figure 4B). *G. leopoldii* and *G. swidsinskii* did not occur together more often than
260 expected by chance, which is of note as both were previously labeled as subgroup A based
261 on cpn60 sequences and whole genome sequence comparisons [8]. The differentiating
262 features of these two species detected by MALDI-TOF [13] may be associated with their
263 occupation of distinct niches. This suggests that the new labeling is indeed useful for
264 understanding differences in distributions at this deeper level. None of the species pairs
265 had fewer co-occurrences than expected, suggesting that competitive exclusion may not be
266 important for describing their relative distributions. Conclusions regarding the rarely
267 detected genome species are limited since very few samples were positive and thus chances
268 of observing co-occurrence were correspondingly low.

269 **Association of *Gardnerella* spp. with BV status and vaginal symptoms**

270 To understand how resolution of different *Gardnerella* spp. may inform clinically
271 important outcomes, we compared the relative abundances of the more frequently
272 occurring species (*G. vaginalis*, *G. swidsinskii*, *G. leopoldii*, *G. piotii* and genome sp. 3)
273 among groups based on clinical Nugent scores (Negative, Intermediate, BV), and self-
274 reported symptoms in the two weeks prior to the swab collection (odour, irritation, and
275 discharge). There was a significant association between Nugent category and relative
276 abundance of *G. vaginalis*, *G. swidsinskii*, and *G. piotii*, (Table 1). Genome sp. 3 was
277 marginally associated with Nugent category, but none of the pairwise comparisons was
278 significant after p-value adjustment. The relationship between *Gardnerella* abundance and
279 Nugent score is not surprising, as the presence of *Gardnerella* “morphotypes” on Gram
280 stained slides of vaginal specimens is part of the calculation of the clinical score [3]. The
281 lack of association of *G. leopoldii* with Nugent category is interesting in that this species
282 was previously investigated together with *G. swidsinskii* as cpn60 subgroup A/clade 4
283 which was found to be associated with Nugent category [14, 34, 40]. These species are
284 very closely related phylogenetically (Figure 1), and were only resolved by Vanechoutte
285 et al. [13] by MALDI-TOF and whole genome comparison, so the specific factors
286 responsible for their apparently different relationships with the microbiological definition
287 of BV remain to be identified.

288 Phenotypic diversity has long been considered a possible explanation for the
289 detection of *Gardnerella* in women without vaginal symptoms, however, attempts to
290 identify associations between particular biotypes and clinical status have yielded
291 inconsistent and often contradictory results [41-45]. The major limitation of investigations
292 relying on phenotypic characterization of isolates is that they focus only on the most readily

293 culturable isolates from individual specimens (often only one isolate per specimen), which
294 is inadequate since women are usually colonized by multiple species of *Gardnerella*. We
295 observed strong relationships between abnormal odour and discharge with higher relative
296 abundance of *G. vaginalis* and *G. swidsinskii*, but not with the other three species, although
297 there is a marginal relationship between discharge and genome sp. 3 ($p = 0.02$) (Table 1).
298 *G. vaginalis* and *G. swidsinskii* co-occurred more often than expected by chance, and also
299 showed proportionality suggesting that they are correlated in abundance. Therefore, we
300 cannot be sure if it is just one species or both that is associated with vaginal symptoms.
301 Sialidase activity defines *G. piotii* [13] and is also observed for genome sp. 3 isolates [8],
302 however, these species were not associated with discharge nor were they the most strongly
303 associated with a BV diagnosis by Nugent score. This is likely due to the polymicrobial
304 nature of BV, and the fact that many other BV-associated bacteria produce hydrolytic
305 enzymes that may contribute to symptoms [46-48]. The lack of association of sialidase
306 positive *Gardnerella* spp. with symptoms is also consistent with the suggestion that some
307 types of *Gardnerella* may be important for “stage-setting”; establishing an anaerobic
308 environment and initial adhesion to the vaginal epithelium that lead to abundant growth of
309 other BV associated organisms, and the development of multi-species biofilms (reviewed
310 in [49]). In these primary colonizers, hydrolytic enzymes and cholesterol-dependent
311 cytolysin (vaginolysin) may be more important in preparing the microbiome for secondary
312 expansion of populations of BV associated bacteria rather than acting as specific virulence
313 factors affecting the host.

314

315 **Conclusions**

316 Considering *Gardnerella* as a monolithic taxon in vaginal microbiome studies (due
317 to the ubiquitous application of 16S rRNA gene sequencing in microbiome profiling) has
318 limited progress in understanding the link between vaginal microbiota and clinical
319 outcomes, and the development of improved diagnostics for women’s health. Our results
320 provide a clear demonstration of the utility of cpn60 barcode sequencing for rapid, high-
321 throughput determination of *Gardnerella* spp. abundance and distribution in the vaginal
322 microbiome with minimal sequencing effort. This approach will be critical in further
323 investigation of the intriguing association of *G. piotii* (subgroup B/clade 2) with
324 “intermediate” microbiota, which has been observed independently using cpn60 barcode
325 sequencing and clade-specific PCR [14, 33]. It remains to be determined if *Gardnerella*
326 species that do not regularly achieve numerical dominance in the microbiome contribute
327 to establishment and maintenance of dysbiosis. Robust and simple classification of
328 *Gardnerella* isolates based on cpn60 barcode sequences will facilitate further
329 characterization of these organisms within the new taxonomic framework, and lead to
330 identification of phenotypic features of the species that determine their ecological roles in
331 the vaginal microbiome. In the clinical context, assessing longitudinal shifts in *Gardnerella*
332 spp. abundance will also be important to evaluate natural or post treatment changes. In
333 future cohort studies, application of cpn60 barcode sequencing will provide new insight as
334 to whether *Gardnerella* spp. diversity and differential distribution are an explanation for
335 issues such as treatment failure and recurrent vaginal dysbiosis [34, 50], or for the failure
336 of antimicrobial treatment in the prevention of preterm birth despite a strong association
337 between vaginal dysbiosis and preterm delivery [51].
338

339

340 **Declarations**

341 *Ethics approval and consent to participate*

342 Studies from which data was accessed for this sub-study were approved by the University
343 of British Columbia Children's & Women's Research Ethics Board (Certificate numbers
344 H10-02535, H11-00119, and H11-01912).

345

346 *Availability of data and materials*

347 The datasets supporting the results of this article is available in the NCBI repository
348 (BioProject Accessions: PRJNA362575, PRJNA278895, PRJNA528096). R code for the
349 co-occurrence analysis and data table are provided as supplemental information.

350

351 *Competing interests*

352 The authors declare that they have no competing interests.

353

354 *Funding*

355 Financial support was provided by a joint Canadian Institutes of Health Research (CIHR)
356 Emerging Team Grant and a Genome British Columbia (GBC) grant (reference #108030)
357 awarded to the VOGUE Research Group, and by an NSERC Discovery Grant to JEH.

358

359 *Authors' contributions*

360 JEH and AYKA conceived the study, conducted the analysis and wrote the paper. The
361 VOGUE Research Group provided access to data for analysis and edited the paper. All
362 authors read and approved the final manuscript.

363

364 *Acknowledgements*

365 The VOGUE Research Group is Deborah Money, Alan Bocking, Sean Hemmingsen,
366 Janet Hill, Gregor Reid, Tim Dumonceaux, Gregory Gloor, Matthew Links, Kieran
367 O’Doherty, Patrick Tang, Julianne van Schalkwyk and Mark Yudin.

368

369 **References**

- 370 1. Leopold S: **Heretofore undescribed organism isolated from the genitourinary**
371 **system.** *US Armed Forces Med J* 1953, **4(2):263-266.**
- 372 2. Schellenberg JJ, Patterson MH, Hill JE: ***Gardnerella vaginalis* diversity and**
373 **ecology in relation to vaginal symptoms.** *Res Microbiol* 2017,
374 **doi:10.1016/j.resmic.2017.02.011.**
- 375 3. Nugent RP, Krohn MA, Hillier SL: **Reliability of diagnosing bacterial vaginosis**
376 **is improved by a standardized method of gram stain interpretation.** *J Clin*
377 *Microbiol* 1991, **29(2):297-301.**
- 378 4. Piot P, Van Dyck E, Peeters M, Hale J, Totten PA, Holmes KK: **Biotypes of**
379 ***Gardnerella vaginalis*.** *J Clin Microbiol* 1984, **20(4):677-679.**
- 380 5. Benito R, Vazquez JA, Berron S, Fenoll A, Saez-Neito JA: **A modified scheme for**
381 **biotyping *Gardnerella vaginalis*.** *J Med Microbiol* 1986, **21(4):357-359.**
- 382 6. Ingianni A, Petruzzelli S, Morandotti G, Pompei R: **Genotypic differentiation of**
383 ***Gardnerella vaginalis* by amplified ribosomal DNA restriction analysis**
384 **(ARDRA).** *FEMS Immunol Med Microbiol* 1997, **18(1):61-66.**
- 385 7. Hill JE, Goh SH, Money DM, Doyle M, Li A, Crosby WL, Links M, Leung A,
386 Chan D, Hemmingsen SM: **Characterization of vaginal microflora of healthy,**
387 **nonpregnant women by chaperonin-60 sequence-based methods.** *Am J Obstet*
388 *Gynecol* 2005, **193(3 Pt 1):682-692.**
- 389 8. Schellenberg JJ, Paramel Jayaprakash T, Withana Gamage N, Patterson MH,
390 Vanechoutte M, Hill JE: ***Gardnerella vaginalis* subgroups defined by cpn60**

- 391 **sequencing and sialidase activity in isolates from Canada, Belgium and Kenya.**
392 *PLoS ONE* 2016, **11**(1):e0146510.
- 393 9. Links MG, Dumonceaux TJ, Hemmingsen SM, Hill JE: **The chaperonin-60**
394 **universal target is a barcode for bacteria that enables *de novo* assembly of**
395 **metagenomic sequence data.** *PLoS ONE* 2012, **7**(11):e49755.
- 396 10. Paramel Jayaprakash T, Schellenberg JJ, Hill JE: **Resolution and characterization**
397 **of distinct cpn60-based subgroups of *Gardnerella vaginalis* in the vaginal**
398 **microbiota.** *PLoS ONE* 2012, **7**(8):e43009.
- 399 11. Ahmed A, Earl J, Retchless A, Hillier SL, Rabe LK, Cherpes TL, Powell E, Janto
400 B, Eutsey R, Hiller NL *et al*: **Comparative genomic analyses of 17 clinical**
401 **isolates of *Gardnerella vaginalis* provide evidence of multiple genetically**
402 **isolated clades consistent with subspeciation into genovars.** *J Bacteriol* 2012,
403 **194**(15):3922-3937.
- 404 12. Richter M, Rossello-Mora R: **Shifting the genomic gold standard for the**
405 **prokaryotic species definition.** *Proc Nat Acad Sci USA* 2009, **106**(45):19126-
406 19131.
- 407 13. Vaneechoutte M, Guschin A, Van Simaey L, Gansemans Y, Van Nieuwerburgh F,
408 Cools P: **Emended description of *Gardnerella vaginalis* and description of**
409 ***Gardnerella leopoldii* sp. nov., *Gardnerella piotii* sp. nov. and *Gardnerella***
410 ***swidsinskii* sp. nov., with delineation of 13 genomic species within the genus**
411 ***Gardnerella*.** *Int J Syst Evol Microbiol* 2019, **69**(3):679-687.
- 412 14. Albert AY, Chaban B, Wagner EC, Schellenberg JJ, Links MG, van Schalkwyk J,
413 Reid G, Hemmingsen SM, Hill JE, Money D: **A study of the vaginal microbiome**

- 414 **in healthy Canadian women utilizing cpn60-based molecular profiling reveals**
415 **distinct *Gardnerella* subgroup community state types.** *PLoS ONE* 2015,
416 **10(8):e0135620.**
- 417 15. Vancuren SJ, Hill JE: **Update on cpnDB: a reference database of chaperonin**
418 **sequences.** *Database (Oxford)* 2019, **2019.**
- 419 16. Hill JE, Penny SL, Crowell KG, Goh SH, Hemmingsen SM: **cpnDB: a chaperonin**
420 **sequence database.** *Genome Res* 2004, **14(8):1669-1675.**
- 421 17. Felsenstein J: **PHYLIP - phylogeny inference package (version 3.2).** *Cladistics*
422 1989, **5:164-166.**
- 423 18. Caporaso JG, Kuczynski J, Stombaugh J, Bittinger K, Bushman FD, Costello EK,
424 Fierer N, Pena AG, Goodrich JK, Gordon JI *et al*: **QIIME allows analysis of high-**
425 **throughput community sequencing data.** *Nature Meth* 2010, **7(5):335-336.**
- 426 19. Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJ, Holmes SP:
427 **DADA2: High-resolution sample inference from Illumina amplicon data.** *Nat*
428 *Methods* 2016, **13(7):581-583.**
- 429 20. Schellenberg J, Links MG, Hill JE, Dumonceaux TJ, Peters GA, Tyler S, Ball B,
430 Severini A, Plummer FA: **Pyrosequencing of the chaperonin-60 universal target**
431 **as a tool for determining microbial community composition.** *Appl Environ*
432 *Microbiol* 2009, **75(9):2889-2898.**
- 433 21. Janson S, Vegelius J: **Measures of ecological association.** *Oecologia* 1981,
434 **49(3):371-376.**

- 435 22. Mainali KP, Bewick S, Thielen P, Mehoke T, Breitwieser FP, Paudel S, Adhikari
436 A, Wolfe J, Slud EV, Karig D *et al*: **Statistical analysis of co-occurrence patterns**
437 **in microbial presence-absence datasets.** *PLoS ONE* 2017, **12**(11):e0187132.
- 438 23. Veech JA: **A probabilistic model for analysing species co-occurrence.** *Global*
439 *Ecology and Biogeography* 2013, **22**:252–260.
- 440 24. Benjamini Y, Hochberg Y: **Controlling the false discovery rate: a practical and**
441 **powerful approach to multiple testing.** *J R Stat Soc Series B Stat Methodol* 1995,
442 **57**:289-300.
- 443 25. Oksanen J, Blanchet FG, Kindt R, Legendre P, Minchin PR, O'Hara RB, Simpson
444 GL, Solymos P, Henry M, Stevens H *et al*: **vegan: Community Ecology Package.**
445 **R package version 2.0-10. Available: [http://CRAN.R-](http://CRAN.R-project.org/package=vegan)**
446 **[project.org/package=vegan](http://CRAN.R-project.org/package=vegan).** 2012.
- 447 26. Fernandes AD, Reid JN, Macklaim JM, McMurrough TA, Edgell DR, Gloor GB:
448 **Unifying the analysis of high-throughput sequencing datasets: characterizing**
449 **RNA-seq, 16S rRNA gene sequencing and selective growth experiments by**
450 **compositional data analysis.** *Microbiome* 2014, **2**:15.
- 451 27. Verbeke TJ, Sparling R, Hill JE, Links MG, Levin D, Dumonceaux TJ: **Predicting**
452 **relatedness of bacterial genomes using the chaperonin-60 universal target**
453 **(cpn60 UT): application to *Thermoanaerobacter* species.** *Syst Appl Microbiol*
454 2011, **34**:171-179.
- 455 28. Katyal I, Chaban B, Hill JE: **Comparative genomics of *cpn60* defined**
456 ***Enterococcus hirae* ecotypes and relationship of gene content differences to**
457 **competitive fitness.** *Microbial Ecol* 2015, **72**(4):917-930.

- 458 29. Chaban B, Links MG, Paramel Jayaprakash T, Wagner EC, Bourque DK, Lohn Z,
459 Albert AYK, van Schalkwyk J, Reid G, Hemmingsen SM *et al*: **Characterization**
460 **of the vaginal microbiota of healthy Canadian women through the menstrual**
461 **cycle.** *Microbiome* 2014, **2**:23.
- 462 30. Freitas AC, Chaban B, Bocking A, Rocco M, Yang S, Hill JE, Money DM: **The**
463 **vaginal microbiome of healthy pregnant women is less rich and diverse with**
464 **lower prevalence of Mollicutes compared to healthy non-pregnant women.** *Sci*
465 *Rep* 2017, **7**:9212.
- 466 31. Freitas AC, Bocking A, Hill JE, Money DM, Group VR: **Increased richness and**
467 **diversity of the vaginal microbiota and spontaneous preterm birth.**
468 *Microbiome* 2018, **6**(1):117.
- 469 32. Schellenberg JJ, Links MG, Hill JE, Dumonceaux TJ, Kimani J, Jaoko W, Wachihi
470 C, Mungai JN, Peters GA, Tyler S *et al*: **Molecular definition of vaginal**
471 **microbiota in East African commercial sex workers.** *Appl Environ Microbiol*
472 2011, **77**(12):4066-4074.
- 473 33. Balashov SV, Mordechai E, Adelson ME, Gyax SE: **Identification,**
474 **quantification and subtyping of *Gardnerella vaginalis* in noncultured clinical**
475 **vaginal samples by quantitative PCR.** *J Med Microbiol* 2014, **63**(Pt 2):162-175.
- 476 34. Hilbert DW, Schuyler JA, Adelson ME, Mordechai E, Sobel JD, Gyax SE:
477 ***Gardnerella vaginalis* population dynamics in bacterial vaginosis.** *Eur J Clin*
478 *Microbiol Infect Dis* 2017.

- 479 35. Gloor GB, Reid G: **Compositional analysis: a valid approach to analyze**
480 **microbiome high-throughput sequencing data.** *Can J Microbiol* 2016,
481 **62(8):692-703.**
- 482 36. Lovell D, Pawlowsky-Glahn V, Egozcue JJ, Marguerat S, Bahler J:
483 **Proportionality: a valid alternative to correlation for relative data.** *PLoS*
484 *Comput Biol* 2015, **11(3):e1004075.**
- 485 37. Erb I, Notredame C: **How should we measure proportionality on relative gene**
486 **expression data?** *Theory Biosci* 2016, **135(1-2):21-36.**
- 487 38. Ulrich W, Almeida-Neto M, Gotelli NJ: **A consumer's guide to nestedness**
488 **analysis.** *Oikos* 2009, **118:3-17.**
- 489 39. Quinn TP, Richardson MF, Lovell D, Crowley TM: **propr: An R-package for**
490 **identifying proportionally abundant features using compositional data**
491 **analysis.** *Sci Rep* 2017, **7(1):16252.**
- 492 40. Shipitsyna E, Krysanova A, Khayrullina G, Shalepo K, Savicheva A, Guschin A,
493 Unemo M: **Quantitation of all four *Gardnerella vaginalis* clades detects**
494 **abnormal vaginal microbiota characteristic of bacterial vaginosis more**
495 **accurately than putative *G. vaginalis* sialidase A gene count.** *Mol Diagn Ther*
496 2019, **23(1):139-147.**
- 497 41. Numanovic F, Hukic M, Nurkic M, Gegic M, Delibegovic Z, Imamovic A, Pasic
498 **S: Importance of isolation and biotypization of *Gardnerella vaginalis* in**
499 **diagnosis of bacterial vaginosis.** *Bosn J Basic Med Sci* 2008, **8(3):270-276.**
- 500 42. Briselden AM, Hillier SL: **Longitudinal study of the biotypes of *Gardnerella***
501 ***vaginalis*.** *J Clin Microbiol* 1990, **28(12):2761-2764.**

- 502 43. Aroutcheva AA, Simoes JA, Behbakht K, Faro S: ***Gardnerella vaginalis* isolated**
503 **from patients with bacterial vaginosis and from patients with healthy vaginal**
504 **ecosystems.** *Clin Infect Dis* 2001, **33**(7):1022-1027.
- 505 44. Tosun I, Alpay Karaoglu S, Ciftci H, Buruk CK, Aydin F, Kilic AO, Erturk M:
506 **Biotypes and antibiotic resistance patterns of *Gardnerella vaginalis* strains**
507 **isolated from healthy women and women with bacterial vaginosis.** *Mikrobiyol*
508 *Bul* 2007, **41**(1):21-27.
- 509 45. Pleckaityte M, Janulaitiene M, Lasickiene R, Zvirbliene A: **Genetic and**
510 **biochemical diversity of *Gardnerella vaginalis* strains isolated from women**
511 **with bacterial vaginosis.** *FEMS Immunol Med Microbiol* 2012, **65**(1):69-77.
- 512 46. Briselden AM, Moncla BJ, Stevens CE, Hillier SL: **Sialidases (neuraminidases)**
513 **in bacterial vaginosis and bacterial vaginosis-associated microflora.** *J Clin*
514 *Microbiol* 1992, **30**(3):663-666.
- 515 47. Wiggins R, Hicks SJ, Soothill PW, Millar MR, Corfield AP: **Mucinases and**
516 **sialidases: their role in the pathogenesis of sexually transmitted infections in**
517 **the female genital tract.** *Sex Transm Infect* 2001, **77**(6):402-408.
- 518 48. Robertson AM, Wiggins R, Horner PJ, Greenwood R, Crowley T, Fernandes A,
519 Berry M, Corfield AP: **A novel bacterial mucinase, glycosulfatase, is associated**
520 **with bacterial vaginosis.** *J Clin Microbiol* 2005, **43**(11):5504-5508.
- 521 49. Hardy L, Cerca N, Jaspers V, Vaneechoutte M, Crucitti T: **Bacterial biofilms in**
522 **the vagina.** *Res Microbiol* 2017, **168**(9-10):865-874.
- 523 50. Bradshaw CS, Morton AN, Hocking J, Garland SM, Morris MB, Moss LM,
524 Horvath LB, Kuzevska I, Fairley CK: **High recurrence rates of bacterial**

525 **vaginosis over the course of 12 months after oral metronidazole therapy and**
526 **factors associated with recurrence. *J Infect Dis* 2006, **193**(11):1478-1486.**

527 51. Nygren P, Fu R, Freeman M, Bougatsos C, Klebanoff M, Guise JM, Force USPST:
528 **Evidence on the benefits and harms of screening and treating pregnant women**
529 **who are asymptomatic for bacterial vaginosis: an update review for the U.S.**
530 **Preventive Services Task Force. *Ann Intern Med* 2008, **148**(3):220-233.**

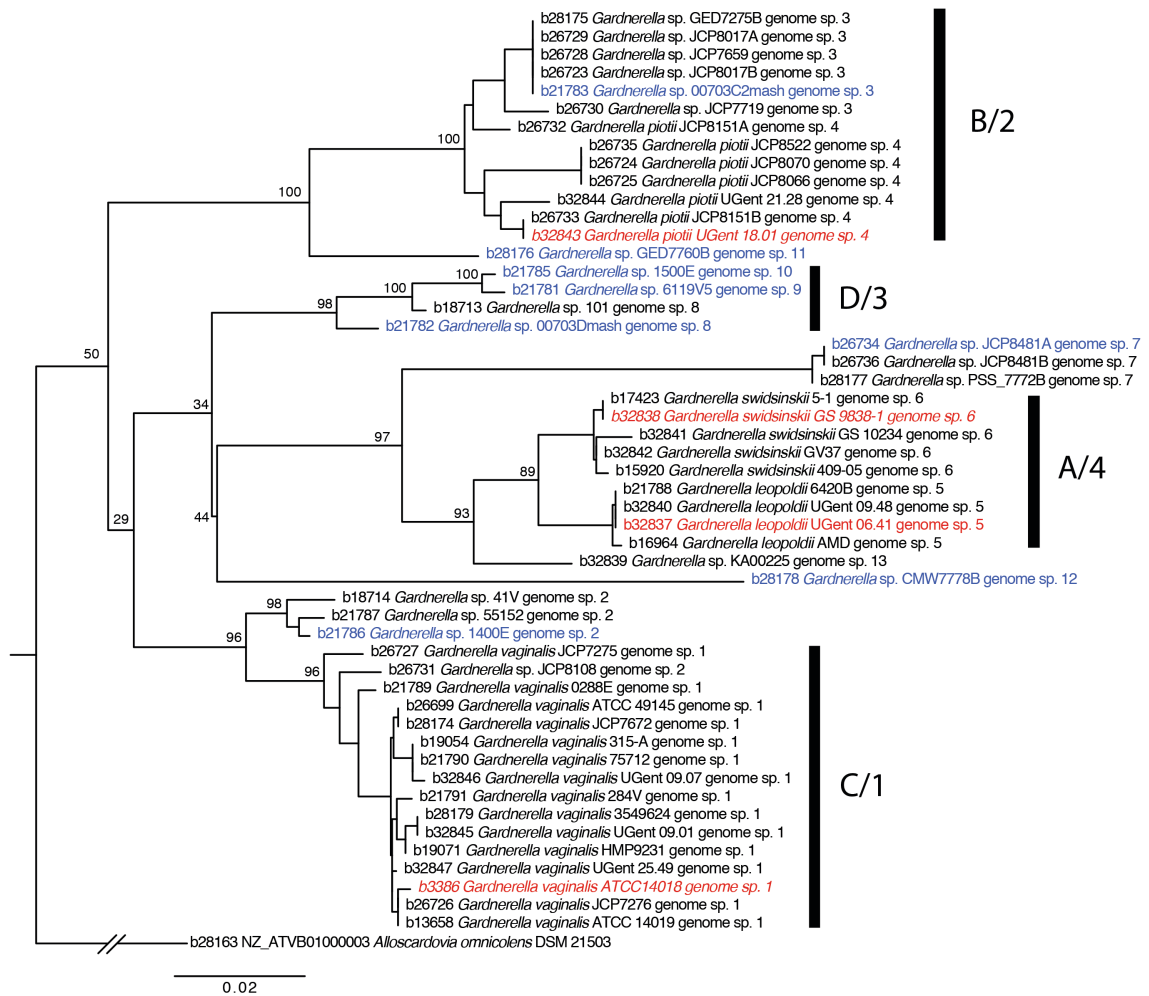
531

532 **Table 1.** Centre-log-ratio (CLR) transformed relative abundance by clinical variables. P-values are from Kruskal-Wallis tests.
 533 Pairwise comparison significance for Nugent category is indicated using letters (a, b) from Dunn tests with Benjamini-Hochberg p-
 534 value adjustment. CLR transform was relative to the geometric mean log₂ relative abundance of all taxa (including non-*Gardnerella*)
 535 therefore negative values indicate relative abundance less than the mean log₂ relative abundance of all taxa.

Variable	<i>G. vaginalis</i>		<i>G. swidsinskii</i>		<i>G. leopoldii</i>		<i>G. piotii</i>		genome sp. 3	
	Median (IQR)	p-value	Median (IQR)	p-value	Median (IQR)	p-value	Median (IQR)	p-value	Median (IQR)	p-value
Nugent		<0.0001		0.0006		0.29		<0.0001		0.03
Negative, n = 289	0.21 (-0.13 to 6.36)	a	0.09 (-0.26 to 5.22)	a	-0.06 (-0.29 to 0.25)	a	-0.06 (-0.31 to 0.14)	a	-0.07 (-0.28 to 0.23)	a
Intermediate, n = 42	7.6 (0.15 to 10.01)	b	0.11 (-0.29 to 8.63)	a,b	-0.06 (-0.24 to 0.44)	a	0.21 (-0.14 to 7.72)	b	0.07 (-0.22 to 4.62)	a
BV, n = 64	9.22 (6.0 to 11.77)	b	4.79 (-0.10 to 12.84)	b	0.002 (-0.28 to 7.64)	a	0.16 (-0.16 to 6.32)	b	0.04 (-0.26 to 6.61)	a
Odour		0.001		0.001		0.11		0.22		0.15
No, n = 360	0.35 (-0.09 to 4.77)		0.10 (-0.25 to 5.89)		-0.06 (-0.29 to 0.28)		-0.02 (-0.29 to 0.24)		-0.06 (-0.27 to 0.31)	
Yes, n = 33	7.56 (4.77 to 10.77)		5.14 (0.03 to 12.80)		0.12 (-0.23 to 8.94)		0.05 (-0.20 to 4.99)		0.03 (-0.24 to 7.56)	
Irritation		0.13		0.20		0.97		0.58		0.25
No, n = 330	0.46 (-0.10 to 8.66)		0.12 (-0.25 to 6.26)		-0.05 (-0.30 to 0.30)		-0.02 (-0.29 to 0.25)		-0.07 (-0.27 to 0.34)	
Yes, n = 63	5.20 (0.02 to 9.11)		0.20 (-0.14 to 8.75)		-0.10 (-0.23 to 0.28)		-0.03 (-0.21 to 0.29)		0.006 (-0.25 to 5.10)	
Discharge		<0.0001		0.001		0.79		0.71		0.02
No, n = 331	0.27 (-0.11 to 8.12)		0.09 (-0.25 to 5.65)		-0.05 (-0.28 to 0.30)		-0.02 (-0.28 to 0.26)		-0.06 (-0.28 to 0.28)	
Yes, n = 62	7.52 (1.55 to 9.83)		2.44 (-0.05 to 11.98)		-0.10 (-0.29 to 0.28)		-0.05 (-0.23 to 0.35)		0.02 (-0.22 to 6.58)	

536
537

538

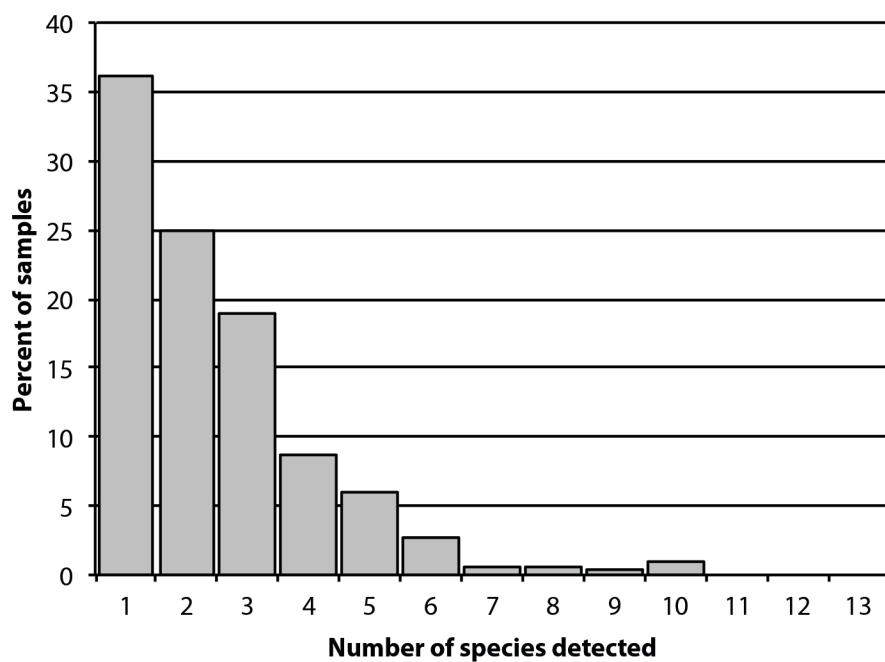


539

540

541 **Figure 1.** Phylogenetic relationships of *Gardnerella* spp. based on an alignment of the
 542 552 bp *cpn60* barcode sequence. Type strains are shown in red, and representatives of
 543 the other 9 genome species designated by Vaneechoutte et al. [13] are in blue. Bootstrap
 544 values are indicated at the major nodes. The tree is rooted with *Alloscardovia*
 545 *omnicolens*. *cpn60* subgroups A-D [10], and clades 1-4 [11] are labeled as
 546 subgroup/clade based on sequences common to previous studies.
 547

548



549

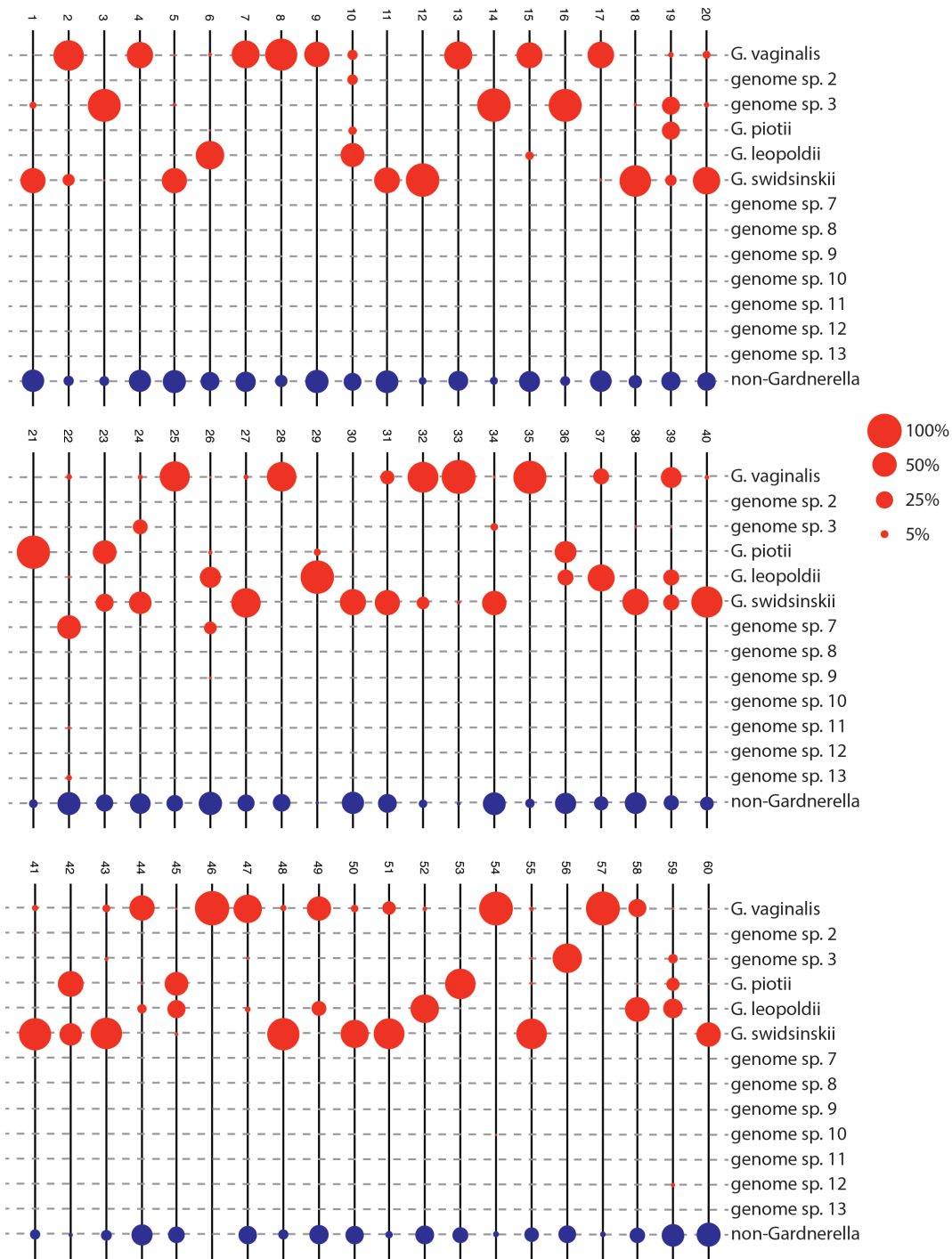
550

Figure 2. Number of *Gardnerella* spp. detected per sample (n = 301).

551

552

553

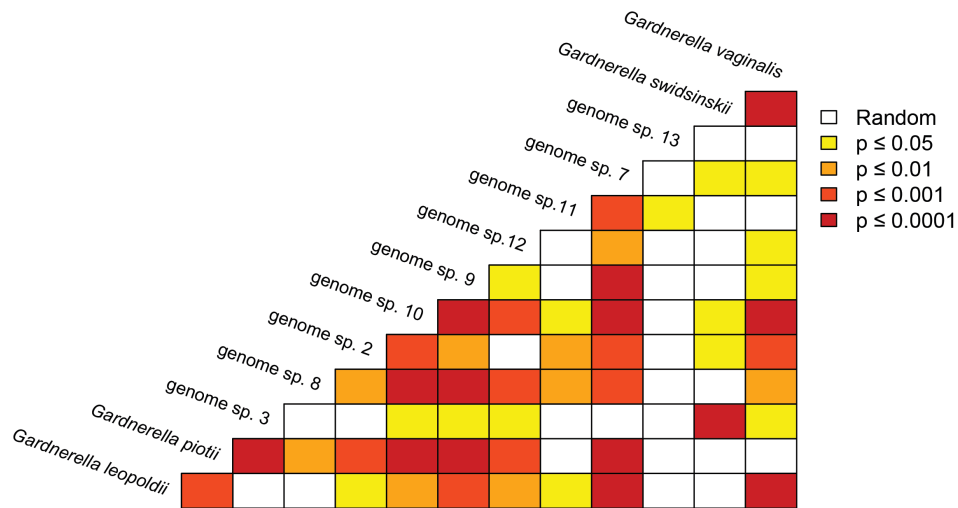


554

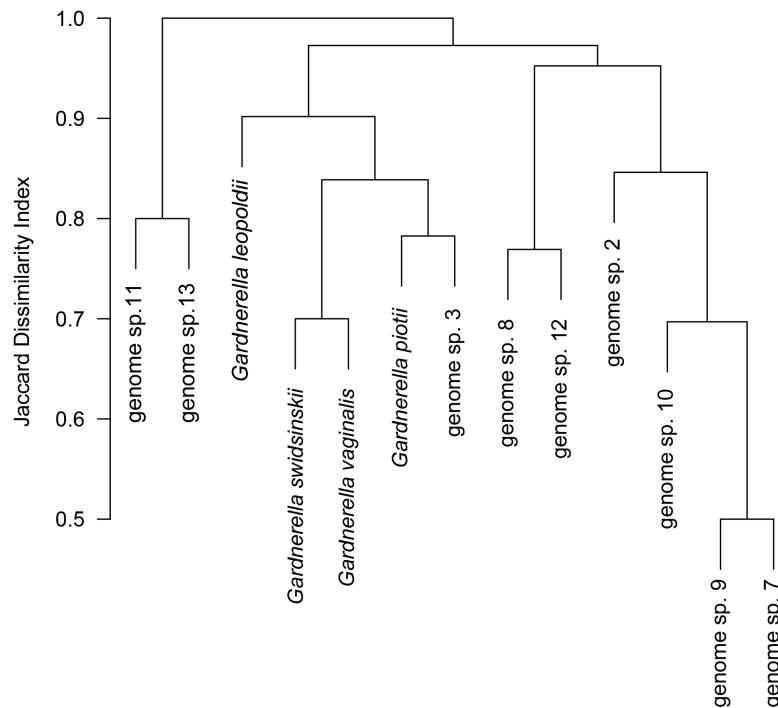
555 **Figure 3.** Proportional abundance of 13 *Gardnerella* spp. in vaginal microbiomes of
556 women with $\geq 50\%$ *Gardnerella* sequence reads (n = 60). Red circles indicate
557 proportional abundance of each species according to scale on the left; blue circles
558 represent the proportion of reads identified as non-*Gardnerella*.
559

560

A.



B.



561

562

563

564

565

566

567

568

Figure 4. (A) Significance of pairwise co-occurrences of *Gardnerella* spp. in 413 vaginal samples determined by Jaccard index of similarity (J) calculation. Size of the P value is indicated by colour according to the legend. Species were considered present if the raw number of sequence reads in the sample was ≥ 10 , otherwise the species was marked as absent. P -values were Benjamini-Hochberg corrected to a false-discovery rate = 0.05 (B) Hierarchical clustering of species based on Jaccard distances ($1 - J$), using complete linkage.

569

570 **Supplemental Information**

571

572 **Table S1.** Pairwise DNA sequence identities among *Gardnerella* spp. based on a
573 CLUSTALw alignment of the 552 bp cpn60 barcode sequence.

574

575 **Table S2.** Sequence read frequencies for 13 *Gardnerella* species in 413 vaginal
576 microbiome samples.

577

578 **Table S3.** Numbers of occurrences and co-occurrences, expected co-occurrences and
579 *P* values for pairwise comparisons of 13 *Gardnerella* species in 413 vaginal samples.

580

581 **Supplemental File 1.** R code for co-occurrence analysis.

582