

1 **Target Journal: Microbiome**

2

3 **High-resolution functional description of vaginal microbiomes in health and disease**

4

5 **Johanna B. Holm^{1,2}, Michael T. France^{1,2}, Pawel Gajer¹, Bing Ma^{1,2}, Rebecca M. Brotman^{1,3}, Michelle**

6 **Shardell^{1,3}, Larry Forney⁴, and Jacques Ravel^{1,2,*}**

7 ¹ Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD, USA

8 ² Department of Microbiology and Immunology, University of Maryland School of Medicine, Baltimore, MD, USA

9 ³ Department of Epidemiology and Public Health, University of Maryland School of Medicine, Baltimore, MD,

10 USA

11 ⁴ Department of Biological Sciences, University of Idaho, Moscow, ID

12 ***Correspondance:**

13 Jacques Ravel

14 jrael@som.umaryland.edu

15

16 **Keywords (3-10 words):**

17 **vaginal microbiome, genital health, metagenome, sequencing, bacterial vaginosis**

18

19 **ABSTRACT**

20 **Background:** A *Lactobacillus*-dominated vaginal microbiome provides the first line of defense against numerous
21 adverse genital tract health outcomes. However, there is limited understanding of the mechanisms by which the
22 vaginal microbiome modulates protection, as prior work mostly described its composition through morphologic
23 assessment and marker gene sequencing methods that do not capture functional information. To address this
24 limitation, we developed metagenomic community state types (mgCSTs) which uses metagenomic sequences to
25 describe and define vaginal microbiomes based on both composition and function.

26

27 **Results:** MgCSTs are categories of microbiomes classified using taxonomy and the functional potential encoded in
28 their metagenomes. MgCSTs reflect unique combinations of metagenomic subspecies (mgSs), which are
29 assemblages of bacterial strains of the same species, within a microbiome. We demonstrate that mgCSTs are
30 associated with demographics such as age and race, as well as vaginal pH and Gram stain assessment of vaginal
31 smears. Importantly, these associations varied between mgCSTs predominated by the same bacterial species. A
32 subset of mgCSTs, including three of the six predominated by *Gardnerella* mgSs, as well as a mgSs of *L. iners*,
33 were associated with a greater likelihood of Amsel bacterial vaginosis diagnosis. This *L. iners* mgSs, among other
34 functional features, encoded enhanced genetic capabilities for epithelial cell attachment that could facilitate
35 cytotoxin-mediated cell lysis. Finally, we report a mgSs and mgCST classifier as an easily applied, standardized
36 method for use by the microbiome research community.

37

38 **Conclusions:**

39 MgCSTs are a novel and easily implemented approach to reducing the dimension of complex metagenomic datasets,
40 while maintaining their functional uniqueness. MgCSTs enable investigation of multiple strains of the same species
41 and the functional diversity in that species. Future investigations of functional diversity may be key to unraveling
42 the pathways by which the vaginal microbiome modulates protection to the genital tract. Importantly, our findings
43 support the hypothesis that functional differences between vaginal microbiomes, including those that may look
44 compositionally similar, are critical considerations in vaginal health. Ultimately, mgCSTs may lead to novel
45 hypotheses concerning the role of the vaginal microbiome in promoting health and disease, and identify targets for
46 novel prognostic, diagnostic, and therapeutic strategies to improve women's genital health.

47

48 **BACKGROUND**

49 The vaginal microbiome plays a vital role in gynecological and reproductive health. *Lactobacillus* predominated
50 vaginal microbiota constitute the first line of defense against infection. Protective mechanisms include lactic acid
51 production by *Lactobacillus* spp., which acidifies the vaginal microenvironment and elicits anti-inflammatory
52 effects [1-4]. This environment wards off non-indigenous organisms, including causative agents of sexually
53 transmitted infections (STIs) like HIV and pathogenic bacteria associated with bacterial vaginosis (BV) [5-7].
54 However, vaginal *Lactobacillus* spp. are functionally diverse. For example, *L. crispatus*, and *L. gasseri* are capable
55 of producing both the D- and L-isomers of lactic acid, *L. jensenii* produces only the D-isomer, and *L. iners* only the
56 L-isomer [4, 8]. These key features have implications for susceptibilities to pathogens [9, 10].

57

58 The vaginal microbiota has been previously shown to cluster into community state types (CSTs) that reflect
59 differences in bacterial species composition and abundance [1, 11]. *Lactobacillus* spp. predominate four of the five
60 CSTs (CST I: *L. crispatus*; CST II: *L. gasseri*; CST III: *L. iners*, CST V: *L. jensenii*). In contrast, CST IV
61 communities are characterized by a paucity of lactobacilli and the presence of a diverse array of anaerobes such as
62 *Gardnerella vaginalis* and “*Ca. Lachnocurva vaginae*”. CST IV is found, albeit not exclusively, during episodes of
63 BV, a condition associated with increased risk to sexually transmitted infections, including HIV, as well as preterm
64 birth and other gynecological and obstetric adverse outcomes [12-20]. BV is clinically defined by observing 3 of 4
65 Amsel’s criteria (Amsel-BV; vaginal pH > 4.5, abnormal discharge, and on wet mount, presence of clue cells and
66 fishy odor with 10% KOH) [21]. Patients presenting with symptoms and satisfying the Amsel’s criteria
67 (symptomatic Amsel-BV) are treated with antibiotics, however, efficacy is poor, and recurrence is common [21-24].
68 In research settings, BV is often defined by scoring Gram stained vaginal smears (Nugent-BV) [25] or molecular
69 typing of bacterial composition by sequencing marker genes (molecular-BV) [26]. There is no definition of BV that
70 relies on both the composition and function of the microbiome.

71

72 Species-level composition of the vaginal microbiota may not suffice to accurately capture associations between the
73 vaginal microbiome and outcomes of interest because functional differences exist between strains of the same
74 species. For example, in the skin microbiome, strains of *Staphylococcus aureus* or *Streptococcus pyogenes* elicit

75 different acute immune responses [27]. Similarly, genomic and functional analyses of *Lactobacillus rhamnosus*
76 strains demonstrate distinct adaptations to specific niches (for example, the gut versus the oral cavity) [28]. While
77 functional differences likely exist between strains of the same species in the vaginal microbiota, metagenomic
78 studies show that combinations of multiple strains co-exist within a single vaginal microbiome [29, 30]. These strain
79 assemblages are known as metagenomic subspecies or mgSs [29], and are important to consider as they potentially
80 impact the functional diversity and resilience of a species in a microbiome. Determining the mechanistic
81 consequences and health outcomes associated with metagenomic subspecies may improve precision of risk estimates
82 and interventions.

83

84 To integrate the taxonomic composition and functional potential of vaginal microbiomes, we developed
85 metagenomic community state types (mgCSTs). MgCSTs are composed of unique combinations of mgSs. We
86 developed and validated a two-step classifier that assigns metagenomic subspecies and mgCSTs and is designed to
87 work in concert with the vaginal non-redundant gene database, VIRGO [29]. This easy-to-use classifier will
88 facilitate reproducibility and comparisons across studies.

89

90 **RESULTS**

91 **Metagenomic community state types (mgCST) of the vaginal microbiome**

92 We evaluated the within-species bacterial genomic diversity in 1,890 vaginal metagenomes of reproductive-age
93 participants from 1,024 mostly North American women (98.7% of samples) (**Table 1 SUBJECT**
94 **DEMOGRAPHICS**). Vaginal metagenomes derived from five cohort studies as well as metagenomes generated to
95 build the vaginal non-redundant gene database (VIRGO, [29]) were used to construct mgCSTs (see **Methods**). In
96 total, 135 metagenomic subspecies (mgSs) from 28 species were identified by hierarchical clustering of species-
97 specific gene presence/absence profiles (**Table S1 Subspecies**). Subsequent clustering of samples based on mgSs
98 compositional data produced 27 mgCSTs (**Table 2 mgCST**). MgCSTs consisted of mgSs from commonly observed
99 vaginal species including *L. crispatus* (mgCST 1-6, 19% of samples), *L. gasseri* (mgCST 7-9, 3% of samples), *L.*
100 *iners* (mgCST 10-14, 23% of samples), *L. jensenii* (mgCST 15 and 16, 4.6% of samples), “*Ca. Lachnocurva*
101 *vaginae*” (mgCST 17-19, 7.5% of samples), *Gardnerella* (mgCST 20-25, 36.3% of samples) and *Bifidobacterium*
102 *breve* (mgCST 26, 0.74% of samples) (**Figure 1 mgCST Heatmap**). MgCST 27 (5.5% of samples) contained less-

103 common species such as *Streptococcus anginosus* or had no predominant taxon. MgCST 2 (n=39 samples from 26
 104 women), mgCST 14 (n=34 samples from 25 women), and mgCST 21 (n=37 samples from 21 women), were only
 105 comprised of samples from reproductive aged women in Alabama enrolled in the UMB-HMP cohort (**Table 2**
 106 **mgCST**). Metagenomic CSTs expand amplicon-based CSTs as multiple mgCSTs are predominated by the same
 107 species, but a different mgSs of that species (**Supplemental Figure 1 Valencia, TABLE 2 mgCST**).
 108

Table 1. Demographic information for all women included in this study. Some women contributed multiple samples.

	Number of Women	Percentage of Women	Number of Samples	Percentage of Samples
Metagenomic Data Source	1,017		1,890	
UMB-HMP	124	12.2	515	27.2
Li <i>et al.</i>	44	35.5	44	8.5
LSVF	585	1329.5	653	1484.1
NIH-HMP	76	13.0	174	26.6
VMRC	40	52.6	162	93.1
VIRGO	148	370.0	342	211.1
Age Category	897		1,623	
15-20	283	31.5	410	25.3
21-25	229	25.5	436	26.9
26-30	188	21.0	362	22.3
31-35	102	11.4	223	13.7
36-40	65	7.2	125	7.7
41-45	30	3.3	67	4.1
Race	858		1,441	
Asian	54	6.3	66	4.6
Black or African American	610	71.1	968	67.2
Hispanic or Latino	19	2.2	47	3.3
Other	6	0.7	9	0.6
White or Caucasian	169	19.7	351	24.4
Nugent Category	968		1,623	
0-3	469	48.5	931	57.4
4-6	194	20.0	255	15.7
7-10	305	31.5	437	26.9
Vaginal pH Category	874		1,362	
Low (pH < 4.5)	273	31.2	491	36.0
High (pH ≥ 4.5)	601	68.8	871	64.0
Amsel-BV Diagnosis	627		673	
Positive	289	46.1	308	45.8
Negative	338	53.9	365	54.2
Symptomatic Amsel-BV	289		308	
Asymptomatic	253	87.5	271	88.0
Symptomatic	36	12.5	37	12.0

Vaginal Non-redundant Gene Database (VIRGO, virgo.igs.umaryland.edu) [29], the University of Maryland Baltimore Human Microbiome Project (UMB-HMP, PRJNA208535, PRJNA575586, PRJNA797778), the National Institutes of Health Human Microbiome Project (NIH-HMP, phs000228), Li *et al.* [60] (PRJEB24147), the Longitudinal Study of Vaginal Flora and Incident STI (LSVF, dbGaP project phs002367).

109



Figure 1. Vaginal Metagenomic Community State Types (mgCSTs). Using 1,890 metagenomic samples, 27 mgCSTs were identified: mgCSTs 1-16 are predominated by metagenomic subspecies of *Lactobacillus* spp., mgCSTs 17-19 by metagenomic subspecies of "Ca. Lachnocurva vaginae", mgCSTs 20-25 by metagenomic subspecies of the genus *Gardnerella*, and mgCST 27 contains samples without a predominant metagenomic subspecies.

Table 2. Metagenomic Community State Types (mgCSTs) of the vaginal microbiome are dominated by different metagenomic subspecies.

MgCST	Most Common mgSs	Most Abundant mgSs	Number of Samples	Number of Women	Median Shannon Index	Number of Samples from Metagenomic Data Source					
						UMB-HMP	Li et al.	LSVF	HMP	VIRGO	VMRC
1	<i>Lactobacillus crispatus</i> 1	<i>Lactobacillus crispatus</i> 1	143	79	0.17	20	2	21	63	15	22
2	<i>Lactobacillus crispatus</i> 2	<i>Lactobacillus crispatus</i> 2	39	26	0.47	39	0	0	0	0	0
3	<i>Lactobacillus crispatus</i> 3	<i>Lactobacillus crispatus</i> 3	83	51	0.28	9	2	15	14	22	21
4	<i>Lactobacillus crispatus</i> 4	<i>Lactobacillus crispatus</i> 4	27	12	0.39	1	1	0	0	3	22
5	<i>Lactobacillus crispatus</i> 5	<i>Lactobacillus crispatus</i> 5	37	27	0.11	1	0	19	16	1	0
6	<i>Lactobacillus crispatus</i> 6	<i>Lactobacillus crispatus</i> 6	28	13	0.69	12	0	5	2	9	0
7	<i>Lactobacillus gasseri</i> 1	<i>Lactobacillus gasseri</i> 1	16	8	0.5	0	0	1	8	1	6
8	<i>Lactobacillus gasseri</i> 2	<i>Lactobacillus gasseri</i> 2	29	17	0.73	15	0	8	0	1	5
9	<i>Lactobacillus gasseri</i> 3	<i>Lactobacillus gasseri</i> 3	14	5	0.89	6	0	0	2	0	6
10	<i>Lactobacillus iners</i> 1	<i>Lactobacillus iners</i> 1	113	76	0.7	24	0	40	2	10	37
11	<i>Lactobacillus iners</i> 2	<i>Lactobacillus iners</i> 2	95	79	0.53	11	7	47	0	28	2
12	<i>Lactobacillus iners</i> 3	<i>Lactobacillus iners</i> 3	131	92	0.44	9	10	45	19	42	6
13	<i>Lactobacillus iners</i> 5	<i>Lactobacillus iners</i> 5	45	41	0.57	1	0	29	1	13	1
14	<i>Lactobacillus iners</i> 6	<i>Lactobacillus iners</i> 6	34	25	0.8	34	0	0	0	0	0
15	<i>Lactobacillus jensenii</i> 1	<i>Lactobacillus jensenii</i> 1	44	28	0.77	8	1	3	10	18	4
16	<i>Lactobacillus jensenii</i> 2	<i>Lactobacillus jensenii</i> 2	67	39	0.71	8	0	15	15	13	16
17	"Ca." <i>Lachnocurva vaginae</i> 1	"Ca." <i>Lachnocurva vaginae</i> 1	58	57	1.48	3	0	51	1	3	0
18	"Ca." <i>Lachnocurva vaginae</i> 1	"Ca." <i>Lachnocurva vaginae</i> 1	28	27	1.57	0	0	27	0	1	0
19	"Ca." <i>Lachnocurva vaginae</i> 1	"Ca." <i>Lachnocurva vaginae</i> 1	43	36	1.91	7	0	27	0	9	0
20	<i>Gardnerella vaginalis</i> 1	<i>Gardnerella vaginalis</i> 1	250	171	1.62	90	2	98	2	38	20
21	<i>Gardnerella vaginalis</i> 1	<i>Gardnerella vaginalis</i> 1	37	21	1.97	37	0	0	0	0	0
22	<i>Gardnerella vaginalis</i> 2	<i>Prevotella amnii</i> 4	202	159	1.79	30	0	91	3	67	11
23	<i>Gardnerella vaginalis</i> 3	<i>Gardnerella vaginalis</i> 3	53	42	0.88	18	5	15	2	5	8
24	<i>Gardnerella vaginalis</i> 4	<i>Gardnerella vaginalis</i> 4	145	106	1.21	44	0	64	6	24	7
25	<i>Gardnerella vaginalis</i> 5	<i>Gardnerella vaginalis</i> 5	34	17	0.83	11	1	9	6	2	5
26	<i>Bifidobacterium breve</i>	<i>Bifidobacterium breve</i>	16	11	0.9	5	0	1	0	8	2
27	<i>Bifidobacterium dentium</i>	<i>Enterococcus faecalis</i> 3	87	76	1.78	23	13	22	2	9	18

UMB-HMP: University of Maryland Baltimore - Human Microbiome Project; Li et al. [62]: PRJEB24147; LSVF: Longitudinal Study of the Vaginal Flora; HMP: Human Microbiome Project; VIRGO: virgo.igs.umaryland.edu; VMRC: Vaginal Microbiome Research Consortium

111

112

113 Vaginal mgCSTs and Demographics

114 **Race and Age.** Race information was available for 1,441 samples from 858 women. Most women identified as

115 either Black (71%) or White (20%), and the remainder identified as Asian (6.3%), Hispanic (2.2%), or other (<1%)

116 (Table 1 SUBJECT DEMOGRAPHICS). Age was also reported for 1,623 samples from 897 individuals and

117 ranged from 15-45 years old. After adjusting for between-cohort heterogeneity, certain races and age categories

118 were associated with mgCSTs (Figure 2). The vaginal microbiomes of Black women were more likely to be

119 classified as *Gardnerella* mgCST 22 ($p = 0.0006$) and least likely to be in *L. crispatus* mgCST 1 ($p = 0.005$) as

120 compared with microbiomes for other races (Table S2 STATS SUMMARY). Microbiomes classified as mgCST 6

121 were more likely to be from White women than other races ($p = 0.002$). *L. iners* mgCST 12 was most common

122 among Hispanic women ($p=0.0001$), and *L. iners* mgCSTs 10 and 14 were absent in Asian women (Figure 2c).

123 MgCSTs predominated by “*Ca. Lachnocurva vaginae*” (mgCSTs 17-19) were also not observed in Asian women,

124 consistent with previous reports on that species (Figure 2c) [11]. In mgCST 27, women were less likely to be Black

125 ($p=0.01$) and more likely to be in the oldest age category (41-45, $p = 0.04$) as compared with other mgCSTs.

126

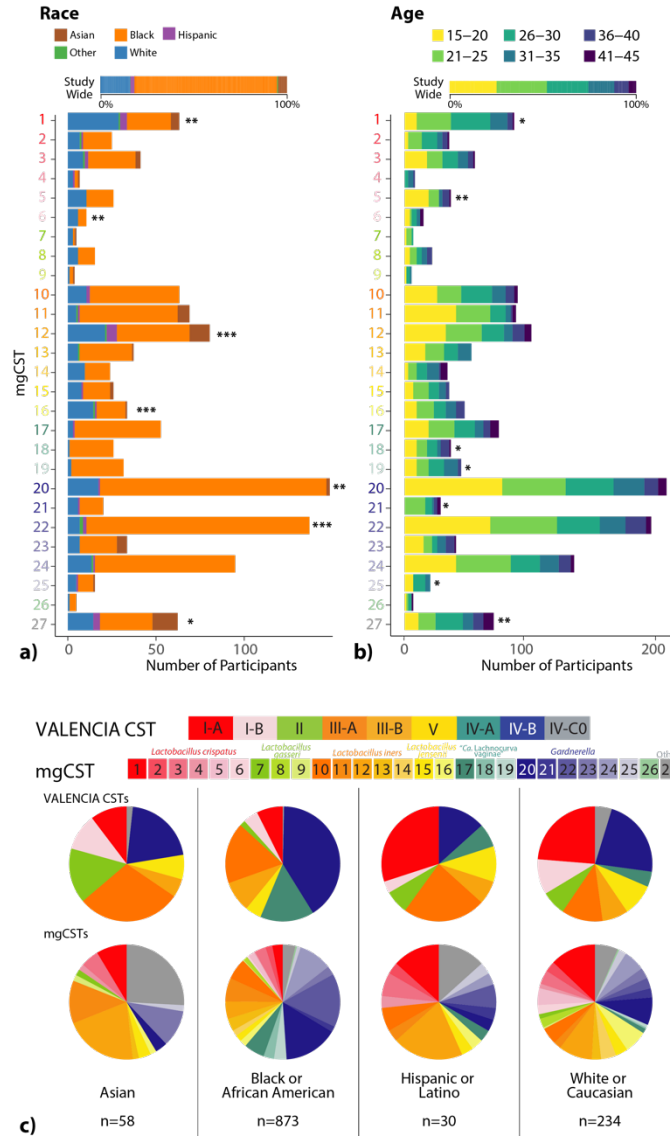


Figure 2. The distribution of (a) race (n=1,441 samples) and (b) age (n=1,623 samples) categories across mgCSTs. Within-mgCST distribution is compared to study-wide distribution (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$). (c) The distribution of mgCSTs across race.

127

128 **Nugent Scores and Vaginal pH.** Of the 968 women for which Nugent scores were available, 48% had low Nugent

129 scores (0-3), 20% had intermediate scores (4-6), and 32% had high scores (7-10) (Table 1). Vaginal pH was also

130 available for 979 women and of these 31% had low pH < 4.5 , and 69% had high pH ≥ 4.5 (Table 1). Both Nugent

131 score and vaginal pH were associated with mgCSTs after adjusting for between-cohort heterogeneity (Figure 3). Of

132 all *L. crispatus* mgCSTs, mgCST 2 had the most representation of different Nugent categories, with 61%, 14%, and

133 25% of samples having low, intermediate, or high Nugent scores, respectively (Figure 3a). Communities

134 predominant in “*Ca. Lachnocurva vaginae*” mgCSTs 17, 18, and 19 had the highest percentages of high Nugent

135 scores (7-10), (94%, 96%, and 87% of samples, respectively); and these mgCSTs were also associated with high
 136 vaginal pH ($p = 6.3 \times 10^{-7}$, **Figure 3b**). Notably, intermediate Nugent scores were common among *Gardnerella*
 137 predominated mgCSTs, especially in mgCSTs 25 (69% of samples).

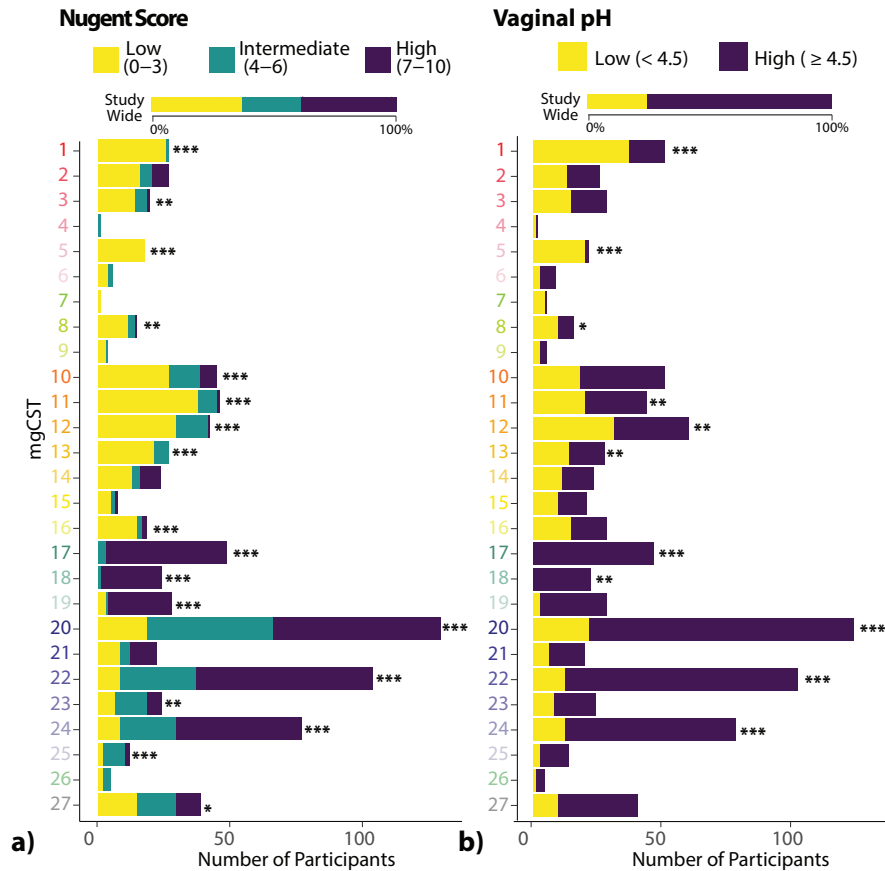


Figure 3. The distribution of (a) Nugent score (n=968), and (b) vaginal pH (n=979) categories. Within-mgCST distribution is compared to study-wide distribution (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$).

138
 139 **Amsel-BV and Vaginal Symptoms.** Of 627 women, each with a sample associated with clinical examination data
 140 (n=607 from LSVF cohort, n=20 from HMP cohort), the proportion of positive Amsel-BV diagnoses (including both
 141 asymptomatic and symptomatic Amsel-BV) was 46%. Twelve percent of Amsel-BV cases were symptomatic.
 142 Diagnosis of Amsel-BV was associated with mgCSTs (**Figure 4a**). There were no Amsel-BV diagnoses in mgCSTs
 143 predominated by *L. crispatus*, *L. jensenii*, or *L. gasseri*. *L. iners* predominated mgCSTs 10-13 were negatively
 144 associated Amsel-BV diagnoses ($p = 9.6 \times 10^{-4}$) but contained some positive Amsel-BV diagnoses in mgCSTs 10, 11,
 145 and 13 (11%, 15%, 18% of women, respectively) (**Figure 4a** and **Table S2 STATS SUMMARY**). *L. iners* mgCST
 146 12 contained only a single (asymptomatic) positive Amsel-BV diagnosis out of 39 women. Women with “Ca.

147 *Lachnocurva vaginae*” mgCSTs 17-19 were more likely to have been diagnosed with Amsel-BV (87%, 88%, and
 148 89%, respectively, $p = 1.8e^{-5}$). *Gardnerella* predominated mgCSTs 20, 22, and 24 also had significantly more
 149 positive Amsel-BV diagnoses than the study-wide proportion (69%, 73%, and 66%, respectively, $p = 1.5e^{-3}$), while
 150 75% of *Gardnerella* predominated mgCST 23 samples were Amsel-BV negative ($p=0.09$). MgCST 24 contained
 151 significantly more symptomatic cases than expected (26% of 43 individuals, $p=0.008$, **Figure 4b, Table S2 STATS**
 152 **SUMMARY**). Though not statistically significant, “*Ca. Lachnocurva vaginae*” mgCST 19 also may have a higher-
 153 than-expected proportion of symptomatic Amsel-BV cases (17.4%).

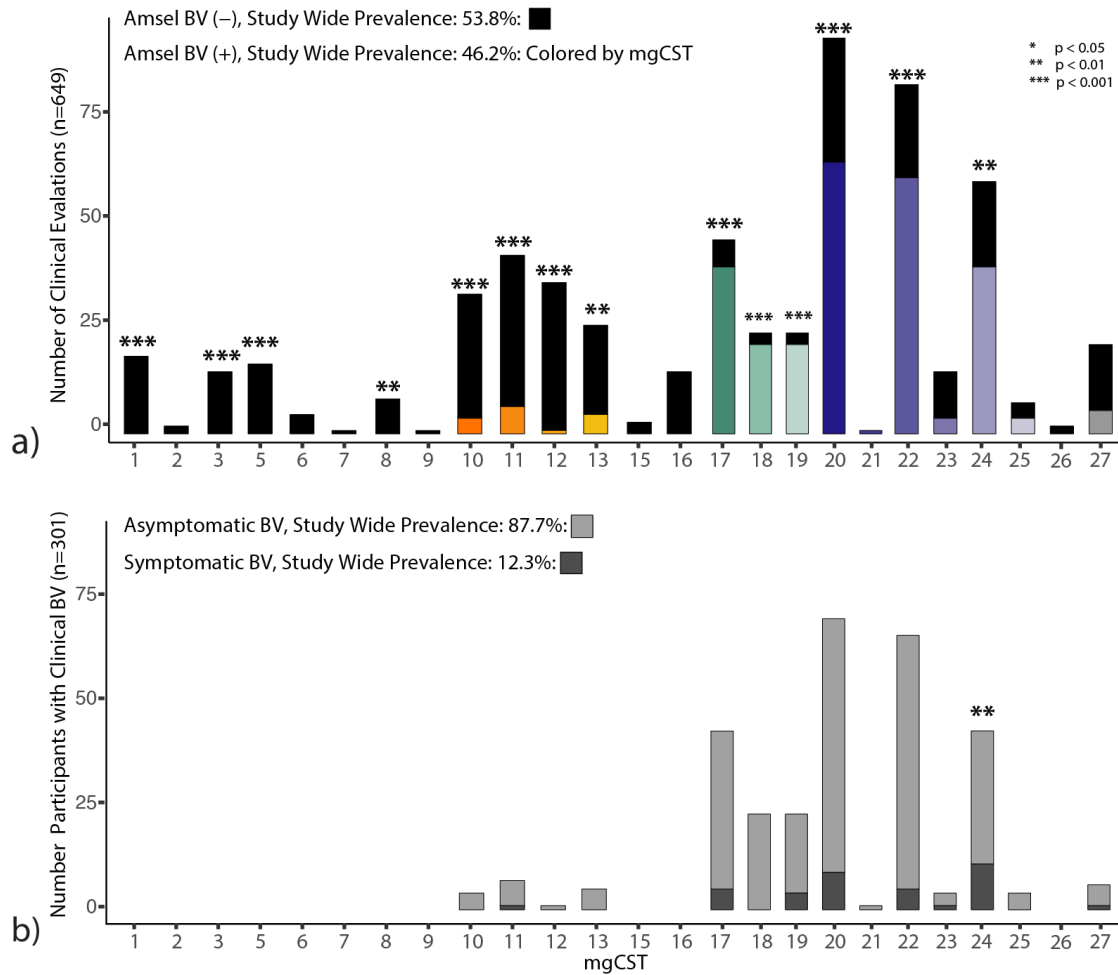


Figure 4. Clinically diagnosed Amsel bacterial vaginosis (a) and symptomatic Amsel bacterial vaginosis (b) are associated with mgCSTs (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$).

154

155

156 **Functional potential of mgCSTs and metagenomic subspecies**

157

158 ***L. crispatus* mgCSTs differ by species diversity, stability, and the potential to produce D-lactic acid.** *L. crispatus* is
159 known to produce both L- and D-lactic acid, which acidifies the vaginal environment and confers protective
160 properties [4, 10, 31, 32]. VIRGO identified two L- and two D-lactate dehydrogenase genes in *L. crispatus*. All
161 genes were present in *L. crispatus* mgCSTs except for mgCST 2. Samples in mgCST 2 were missing a D-lactate
162 dehydrogenase gene (V1806611) that has 96.1% identity to a functionally validated ortholog, P30901.2 (**Figure 5a**)
163 [33]. The other D-lactate dehydrogenase, V1891370, is found in all *L. crispatus* mgCSTs but only 82.4% identical to
164 P30901.2. It contains a 55 aa insertion after V101 (position in P30901.2) and a point mutation at position 218
165 (D218Y) located within a NAD binding site domain. The absence of V1806611 may have functional consequences
166 for microbiomes in mgCST 2. Additionally, samples in mgCST 2 have fewer estimated numbers of *L. crispatus*
167 strains compared to other mgCSTs (**Figure 5b**). Thus, it is likely that an *L. crispatus* strain (or strains) containing
168 V180661 is absent from mgCST 2 samples. Interestingly, the median vaginal pH in mgCST 2 was 4.7, while in
169 mgCST 1 it is 4.0 (1st-3rd quartile: 3.8-4.2, **Figure 5c**). Correspondingly, mgCST 2 samples contained a higher
170 Shannon's H index than mgCST 1 (**Figure 5d**). All samples in mgCST 2 contained genes from "*Ca. Lachnocurva*
171 *vaginae*", *Fingoldia magna*, *Peptoniphilus harei*, *P. lacrimalis*, *Prevotella timonensis*, *P. disiens*, *P. buccalis*, and
172 *Propionibacterium*, albeit at low relative abundances (<1%). We hypothesized that the observed heterogeneity in the
173 compositions of mgCST 2 might result in lower microbiome stability than mgCST 1. Using longitudinal data from
174 the UMB-HMP study, Yue-Clayton θ of daily bacterial composition data over 10 weeks was calculated as an
175 estimate of community stability. Compared to mgCST 1, mgCST 2 samples were indeed significantly less stable ($t =$
176 4.073, $df = 47.942$, $p\text{-value} < 0.001$, **Figure 5e**).

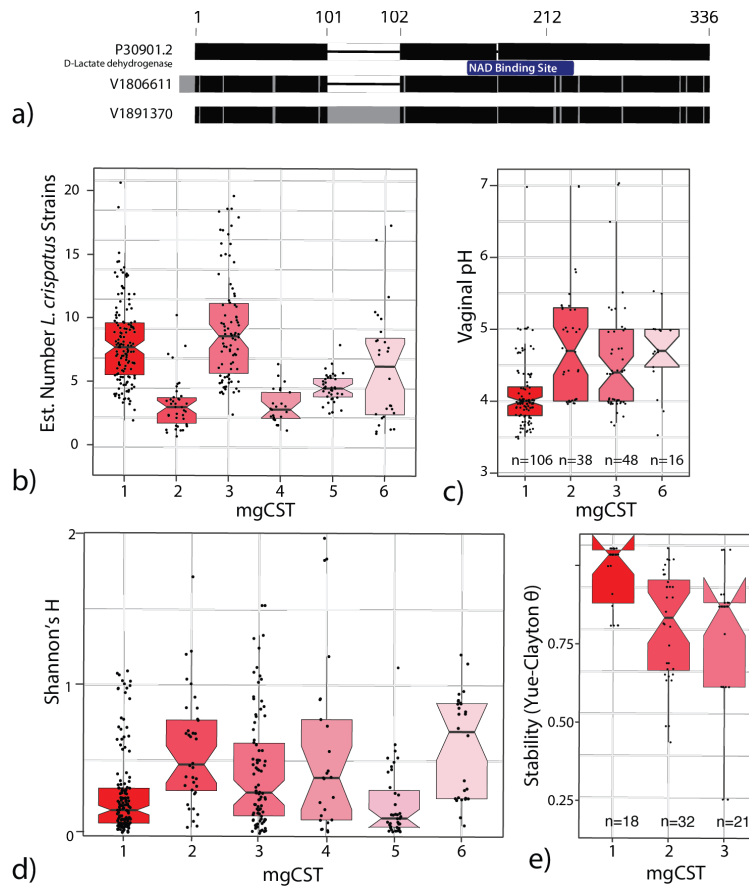


Figure 5. a) D-lactate dehydrogenase orthologs in VIRGO compared to reference P30901.2. b) MgCST 2 contains fewer estimated strains of *L. crispatus*. c) On average, vaginal pH is higher in mgCST 2. d) Shannon's H is higher in mgCST 2 than mgCST 1 or 3. e) Microbiome stability is lower in mgCST 2.

177

178 ***L. iners* metagenomic subspecies are associated with Amsel-BV diagnoses.** The role of *L. iners* in the vaginal

179 microbiome is not fully understood because it is implicated in both healthy and diseased states [34]. *L. iners* is

180 represented by six mgSs. Predominance by *L. iners* mgSs 4 did not define a mgCST (**Figure 1**). Instead, *L. iners* mgSs

181 4 was present in relatively lower abundances (median: 1.2%, IQR: 1.9%) in 257 microbiomes from BV-like mgCSTs

182 including "*Ca. Lachnocurva vaginae*" mgCSTs 16, 17, and 18, and *Gardnerella* mgCSTs 19 and 24. Seventy percent

183 of samples containing *L. iners* mgSs 4 were positive Amsel-BV cases which is significantly greater than the proportion

184 of cases harboring any *L. iners* mgSs (45.8%, $p=1.1 \times 10^{-6}$, **Figure 6a**). Conversely, *L. iners* mgSs 3 was associated with

185 negative Amsel-BV diagnoses (92% Amsel-BV negative, $p=1.6 \times 10^{-9}$).

186

187 We next evaluated if *L. iners* genes were associated with Amsel-BV. Most samples in *L. iners* mgSs 4 contained genes

188 from cluster 6 (yellow gene cluster, **Figure 6b**). There were significantly more positive Amsel-BV diagnoses among

189 subjects containing *L. iners* gene cluster 6 (69.4%, $p=2.1 \times 10^{-15}$), 7 (53.9%, $p=0.004$), or 8 (60.2%, $p=0.036$) compared

190 to samples containing any other *L. iners* gene cluster (45.8%, **Figure 6c**). Gene products unique to *L. iners* gene cluster
 191 6 had significant similarity to virulence factors that could contribute to *L. iners* ability to thrive in dynamic vaginal
 192 states. Such factors include serine/threonine-protein kinases (STPKs), SHIRT domains known as “periscope proteins”
 193 which regulate bacterial cell surface interactions related to host colonization [35], CRISPR-*cas*, β -lactamase and
 194 multidrug resistance (MATE), and bacteriocin exporters (**Table S3**). Gene products in cluster 7 included ParM, which
 195 plays a vital role in plasmid segregation, pre-protein translocation and membrane anchoring (SecA, SecY, sortase),
 196 defense mechanism beta-lytic metalloproteinase, and mucin-binding and internalin proteins. In *Listeria*
 197 *monocytogenes*, internalin A mediates adhesion to epithelial cells and host cell invasion [36]. Phage-like proteins in
 198 gene group 8 suggest the presence of mobile elements. The presence of the highly-conserved *L. iners* pore-forming
 199 cytolysin, inerolysin [37], did not differ by mgSs.

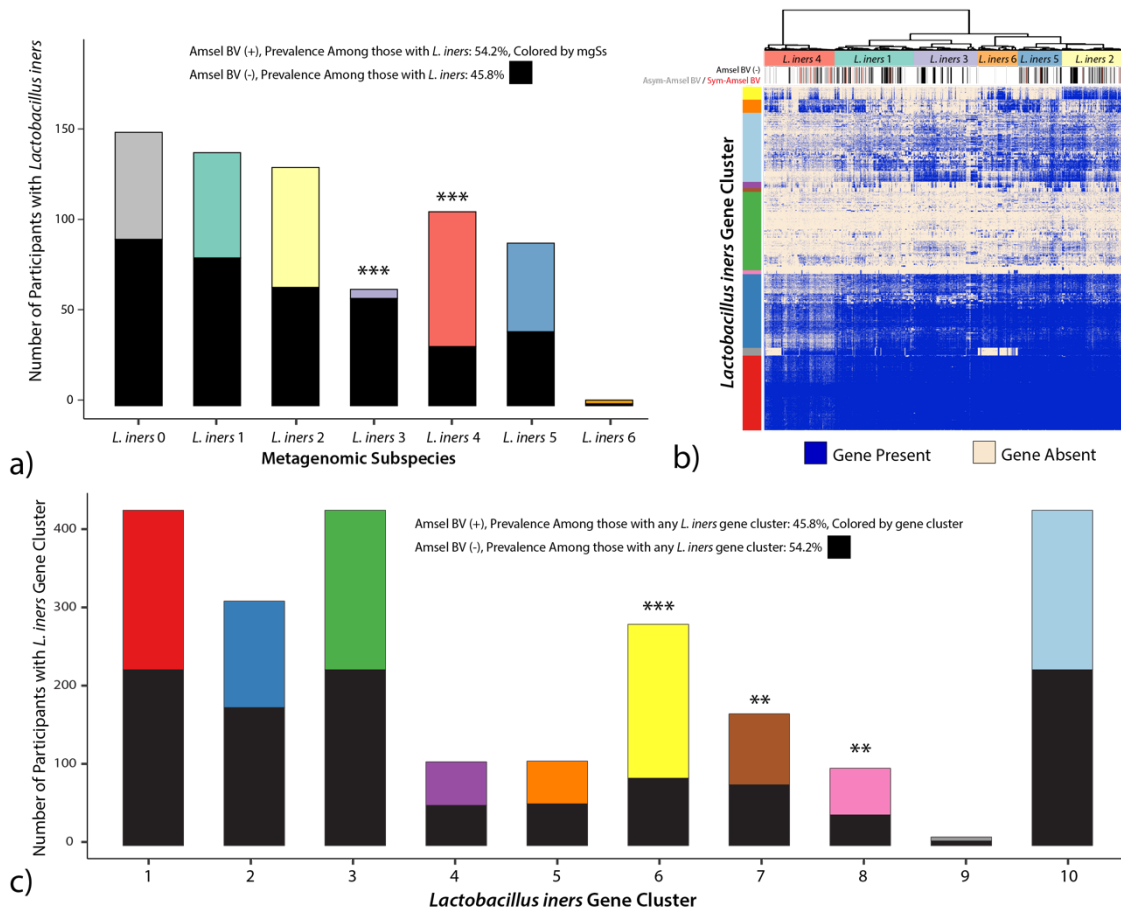


Figure 6. a) Clinically diagnosed Amsel-BV is associated with *L. iners* metagenomic subspecies (mgSs). b) Gene clusters present in *L. iners* mgSs. c) *L. iners* gene clusters 6 (yellow), 7 (brown) and 8 (pink) are associated with positive Amsel-BV diagnosis. Gene cluster 2 (dark blue) is associated with negative Amsel-BV diagnoses. (*: $p < 0.05$, **: $p < 0.01$, ***: $p < 0.001$).

201
 202 **Diversity of Gardnerella genomospecies is associated with an increase in virulence factors.** As previously
 203 mentioned, positive Amsel-BV diagnoses were common in *Gardnerella* mgCSTs 20, 22, and 24, while mgCST 23
 204 contained more negative Amsel-BV diagnoses. Symptomatic BV cases were more common in mgCST 24 than in all
 205 other mgCSTs. By mapping the available genomes of various *Gardnerella* genomospecies [38] to VIRGO, we
 206 determined that each *Gardnerella* mgSs consists of a unique combination of *Gardnerella* genomospecies (**Figure 7a**).
 207 Compared to other *Gardnerella* mgCSTs, mgCSTs 20-22 contain a greater number of *Gardnerella* genomospecies
 208 than mgCSTs 23-25. MgCST 24 samples are predominated by *Gardnerella* mgSs 4 and largely consists of *G.*
 209 *swidsinkii* and *G. vaginalis* genes. This suggests the diversity and types of *Gardnerella* genomospecies may be
 210 important determinants of the pathogenicity of mgCSTs. For example, there are more gene variants of common
 211 *Gardnerella* virulence factors like sialidase and vaginolysin in samples with more genomospecies (**Figure 7b**).

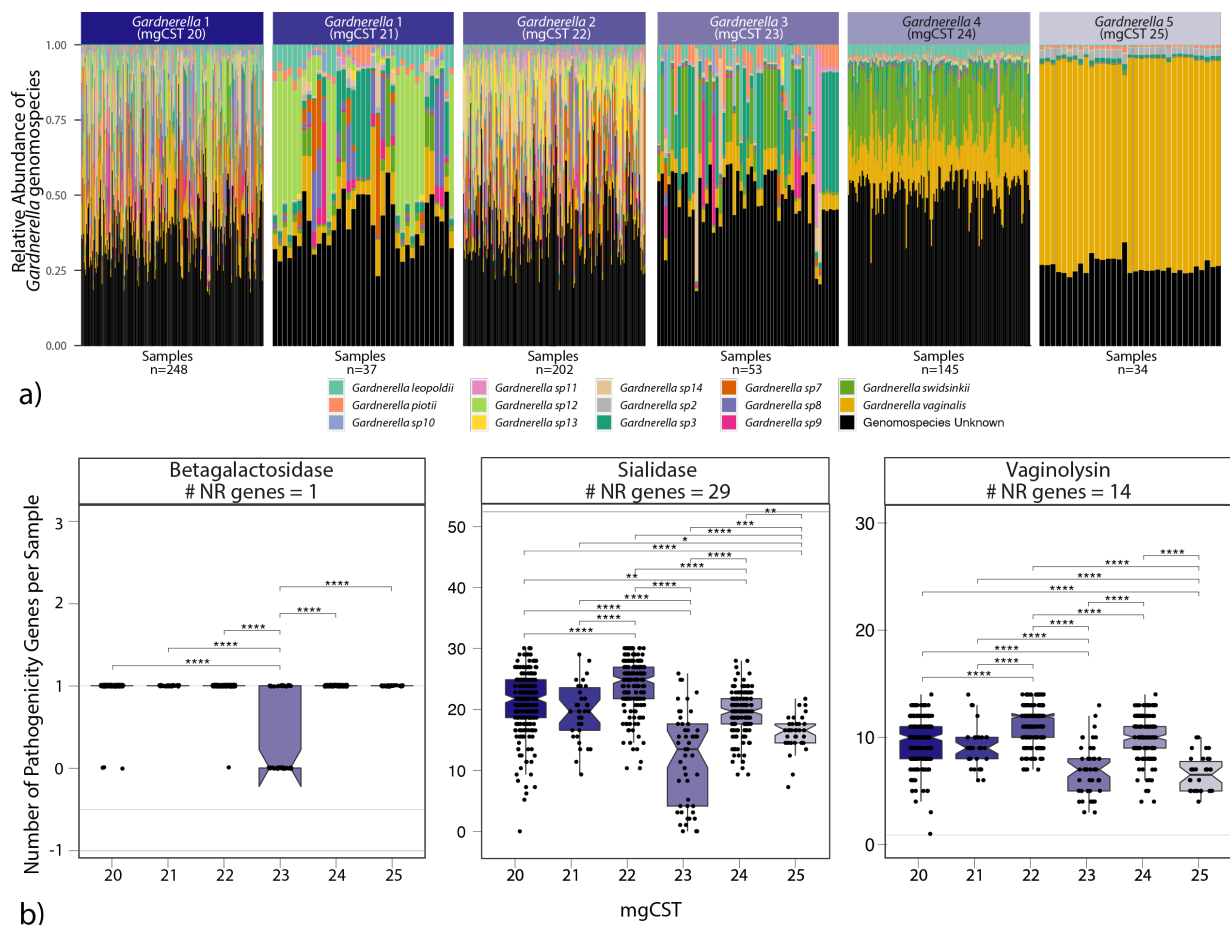


Figure 7. a) The distribution of *Gardnerella* genomospecies across *Gardnerella* mgSs. b) known pathogenicity genes are differently distributed across *Gardnerella* mgCSTs.

212

213 **Automated Classification of mgCSTs using Random Forest Models**

214 Random forest models were built for each of the 135 mgSs identified and used to perform mgSs assignments (**see**
215 **Methods**). The misclassification error for mgSs assignment ranged from 0-30% (**Supplemental Figure 2 mgSs**
216 **Misclassification Error**). The error estimates for most major vaginal taxa were near or less than 10%, with *L.*
217 *gasseri* having the lowest (2.2%). *L. iners* consistently provided higher misclassification error estimates (20%)
218 regardless of attempts to fine-tune the model and was likely the result of high genetic heterogeneity within *L. iners*
219 mgSs. Following assignment of mgSs, mgCSTs were assigned using the nearest centroid classification method, as
220 previously used for vaginal taxonomy-based community state type assignments [11]. The mean classification error
221 was 9.6%, with some mgCSTs classified more accurately than others (**Supplemental Figure 3 mgCST**
222 **Misclassification Error**). The mgCST classifier is packaged into an R script, is available at
223 <https://github.com/ravel-lab/mgCST-classifier> and uses direct outputs from VIRGO.

224 **DISCUSSION**

225 Recent findings that motivated development of mgCST classification are that multiple strains of the same species
226 are commonly observed in the vaginal microbiome [29], and that samples can be clustered into metagenomic
227 subspecies defined by unique strain combinations represented by species-specific gene sets, and thus unique sets of
228 functions. These critical observations led us to conceptualize a vaginal microbiomes classification based on their
229 mgSs compositions and abundance, and thus defined by both species' composition and functions, *i.e.*, metagenomic
230 community state types. MgCSTs describe vaginal microbiomes through a new lens, one that includes both
231 compositional and functional dimensions.

232
233 *L. iners*-predominated vaginal microbiota have been associated with increased risks of experiencing bacterial
234 vaginosis (BV) [39, 40]. Longitudinal observational prospective studies support this conclusion and present several
235 critical findings: 1) *L. iners* is often detected at low to medium abundances during episodes of BV, and *L. iners*
236 commonly dominate the vaginal microbiota after metronidazole treatment for BV and, 2) *L. iners* predominated
237 vaginal microbiota are more prevalent prior to incidence of BV [41, 42]. We observed the frequency of *L. iners*
238 predominated vaginal microbiota was high in Black and Hispanic women (31.4% and 36.1%, respectively), both of
239 whom experience a disproportionate prevalence of BV in the US, with reported rates of 33.2% and 30.7%,
240 respectively (compared to 22.7% and 11.1% in White and Asian women) [43]. Interestingly, *L. iners* predominated

241 vaginal microbiota were even more frequent in North American Asian women in this study, as was shown
242 previously by Ravel *et al.* [1], yet these *L. iners* predominated vaginal microbiota are not associated with higher
243 risk of BV in these women [1]. MgCST classification provides insight into this contradiction to prevailing dogma
244 regarding *L. iners* and increased risk of BV. We noted the absence of *L. iners* mgSs 4 in Asian women, and that *L.*
245 *iners* mgSs 4 is associated with Amsel-BV, while *L. iners* mgSs 3 (predominates mgCST 12) was significantly
246 negatively related to BV and was most frequently observed in Asian women. This is the first evidence of
247 genetically distinct combinations of *L. iners* strains (mgSs) in healthy versus BV-associated states of the vaginal
248 microbiome. This critical finding points to the possibility of beneficial properties associated with some *L. iners*-
249 dominated microbiomes that had not been evidenced previously. Our analyses also identified a specific set of *L.*
250 *iners* genes associated with positive Amsel-BV diagnoses. Macklaim *et al.* 2018 reported marked differences in *L.*
251 *iners* gene expression between two control patients versus two diagnosed with BV, including increased CRISPR-
252 associated proteins gene expression in BV samples [44]. However, our mgSs analysis of *L. iners* indicates that it is
253 not simply alterations in gene expression of a common gene pool that differentiates BV from non-BV microbiomes,
254 but *L. iners* mgSs that also differ. Microbiomes from women with positive BV diagnoses were enriched for host
255 immune response evasion and host-colonization functions. For example, serine/threonine-protein kinases (STPKs)
256 contribute to resistance from phagocytosis by macrophage, invasion of host cells including epithelia and
257 keratinocytes, antibiotic resistance, disruption of the NF- κ B signaling pathway, and mucin binding [45]. Bacterial
258 attachment to host cells (clue cells) is a hallmark of high Nugent scores (a bacterial morphology-based definition of
259 bacterial vaginosis) and a criterion in Amsel-BV diagnoses [21, 25]. Attachment of *L. iners* to epithelial cells may
260 look like clue cells and this could lead to morphological misdiagnosis of BV. These features may make *L. iners*
261 mgSs 4 more difficult to displace in the vaginal environment and could contribute to the common observation of *L.*
262 *iners* following antibiotic treatment [46]. Interestingly, just like *L. iners* mgSs 4, mgSs of “*Ca. Lachnocurva*
263 *vaginae*” were strongly associated with Amsel-BV and were also not found in the vaginal microbiomes of Asian
264 women in this study. Together, these observations may be evidence of selective pressures by the host environment
265 or niche specialization by vaginal bacteria. Sources of selective pressure could relate to host-provided nutrient
266 availability (*e.g.*, mucus glycan composition), the host innate and adaptive immune system, the circulation of other
267 species’ mgSs in a population, or any such combination.

268

269 Several distinct mgCSTs associate strongly with Amsel-BV. Critically, these data support the need for an improved
270 definition of BV and the importance of a personalized approach to treatment. “*Ca. Lachnocurva vaginae*”
271 predominated mgCSTs were strongly associated with asymptomatic Amsel-BV and contained more high Nugent
272 scores than other mgCSTs. Conversely, intermediate Nugent scores were most prevalent in *Gardnerella*
273 predominated mgCSTs, and only three of these six mgCSTs were associated with Amsel-BV, which suggests that
274 not all *Gardnerella*-dominated microbiomes are related to Amsel-BV. *Gardnerella* contains vast genomic diversity,
275 supporting a split into different genomospecies [38, 47, 48]. Because different genomospecies can co-exist, it is
276 likely that *Gardnerella* predominated mgCs represent unique combinations of genomospecies and strains of these
277 genomospecies. Greater *Gardnerella* genomospecies diversity is associated with positive Amsel-BV diagnoses in
278 studies using qPCR or transcriptomic data to define *Gardnerella* species [47, 49-51]. Our data corroborate these
279 reports and further indicate in mgCSTs with higher numbers of *Gardnerella* genomospecies that there are more
280 gene variants coding for virulence factors like cholesterol-dependent pore-forming cytotoxin vaginolysin and
281 neuraminidase sialidase present, thus expanding functional diversity and potentially explaining the association with
282 positive Amsel-BV diagnoses [52-54]. Enumeration of *Gardnerella* genomospecies may prove to be an important
283 diagnostic of certain “types” of Amsel-BV and could inform treatment options. For example, it is possible that
284 harboring more *Gardnerella* genomospecies may predict BV recurrence following metronidazole treatment,
285 suggesting the need for a different approach to treatment. Alternatively, some *Gardnerella* genomospecies may be
286 important and novel targets of therapy.

287
288 In the clinic, antibiotic treatment is recommended for BV diagnosis (generally a point-of-care test) only when the
289 patient reports symptoms, which is estimated to occur in fewer than half of women with BV [24, 55, 56]. In
290 research settings, both symptomatic and asymptomatic Amsel-BV can be evaluated. Indeed, in the observational
291 research studies included in this analysis where Amsel criteria were evaluated along with whether participants
292 reported symptoms or not, symptomatic Amsel-BV accounted for only 12% of Amsel-BV cases and 30% of these
293 were in mgCST 24 (dominated primarily by *Gardnerella swidsinkii* and *G. vaginalis*). We hypothesize that the
294 inadequacy of currently recommended BV treatment may be due to the heterogeneity in the genetic make-up of the
295 microbiota associated with BV as revealed by mgCSTs. MgCSTs reduce this heterogeneity resulting in more

296 precise estimates of risk. Furthermore, these findings highlight the potential importance of developing specialized
297 treatments that target “types” of BV.

298

299 The mgCST framework can also be used to identify vaginal microbiomes that are associated with positive health
300 outcomes. For example, mgCSTs predominated by different *L. crispatus* mgSs varied in their association with low
301 Nugent scores, the number of *L. crispatus* strains present, and the longitudinal stability of communities. The vaginal
302 microbiome can be dynamic [57-59]. Shifts from *Lactobacillus* to non-*Lactobacillus* predominated microbiota can
303 increase the risk of infection following exposure to a pathogen. Our study identified *L. crispatus* mgCSTs with
304 variable stability, suggesting that not all *L. crispatus* predominated microbiomes are functionally similar and may
305 be differently permissive to infection. Those found to be associated with higher stability may reduce the window of
306 opportunity for pathogens to invade. Microbiome stability may be related to both the diversity of other non-
307 *Lactobacillus* members of the microbiome and/or the number of *L. crispatus* strains present. In any case, our study
308 shows that there is a range of protective abilities even among *L. crispatus* predominated communities. This
309 information could be critical in selecting and assembling strains of *L. crispatus* to design novel live biotherapeutics
310 products aimed to restore an optimal vaginal microenvironment.

311

312 It is unclear what factors contribute to vaginal strain assemblages and what rules define their biology and ecology.
313 However, such assemblages can now be detected and further characterized using the concepts of mgSs and
314 mgCSTs presented here. The use of metagenomic sequencing and mgSs and mgCSTs will contribute to a much-
315 needed functional understanding of the role of the vaginal microbiome in reproductive health outcomes. Our
316 findings support the hypothesis that genetic and functional differences between vaginal microbiomes, including
317 those that may look compositionally similar, are critical considerations in vaginal health [7]. To aid in further
318 exploration, we also provide a validated classifier for both mgSs and mgCSTs at [https://github.com/ravel-](https://github.com/ravel-lab/mgCST-classifier/blob/main/README.md)
319 [lab/mgCST-classifier/blob/main/README.md](https://github.com/ravel-lab/mgCST-classifier/blob/main/README.md).

320

321 **CONCLUSION**

322 MgCSTs reveal differences between vaginal microbiome both compositionally and functionally, and thus more
323 finely describe the vaginal microbiome. Associations between mgCSTs and bacterial vaginosis highlight the multi-

324 faceted aspects of the condition and call for new and expanded definitions. Further, we provide tools for the
325 classification of mgSs and mgCST that have potential for use and harmonization of analytical strategies in future
326 studies.

327

328 **DATA AVAILABILITY**

329 The classifiers are available to accompany VIRGO output at <https://github.com/ravel-lab/mgCST-classifier>.

330

331 **COMPETING INTERESTS STATEMENT**

332 JR is co-founder of LUCA Biologics, a biotechnology company focusing on translating microbiome research into
333 live biotherapeutics drugs for women's health. JR is Editor-in-Chief at *Microbiome*. All other authors declare that
334 they have no competing interests.

335

336 **FUNDING**

337 Research reported in this publication was supported in part by the National Institute for Allergy and Infectious
338 Diseases of the National Institutes of Health under award numbers F32-AI136400 (JH), K01-AI163413 (JH),
339 U19AI084044 (JR), UH2AI083264 (JR), R01-AI116799 (RB), and the National Institute for Nursing Research of
340 the National Institutes of Health under award number R01NR015495 (JR). The funders had no role in study design,
341 data collection and interpretation, or the decision to submit the work for publication.

342

343 **METHODS**

344 **Study cohorts.** Raw metagenomic data from 1,890 vaginal samples were used in this study (**Supplemental File 6**).
345 This included publicly available metagenomes including those used in the construction of the vaginal non-redundant
346 gene database, VIRGO (virgo.igs.umaryland.edu, n=342) [29], the University of Maryland Baltimore Human
347 Microbiome Project (UMB-HMP, n=677, PRJNA208535, PRJNA575586, PRJNA797778)[41], the National
348 Institutes of Health Human Microbiome Project (NIH HMP, n=174, phs000228) [60], metagenomes from Li *et al.*
349 [61] (n=44, PRJEB24147), the Longitudinal Study of Vaginal Flora and Incident STI (LSVF, n=653, dbGaP project
350 phs002367) [24]. All samples in LSVF (n=653) and some in UMB-HMP (n=20) had clinical diagnosis information
351 about Amsel-BV. Amsel-BV was diagnosed based on the presence of 3 out of 4 Amsel's criteria [21] and

352 symptomatic Amsel-BV was diagnosed when a woman reported symptoms upon questioning [56]. At the time of
353 these studies, gender identity information was not collected. We know all women responded to recruiting materials
354 which included “women” or “woman”. In addition, individuals are referred to as women in previous publications,
355 thus we refer here to individuals as “woman” or “women” to maintain consistency.

356

357 **Sequence Processing and Bioinformatics.** Host reads were removed from all metagenomic sequencing data using
358 BMTagger and the GRCh38 reference genome, and reads were quality filtered using trimmomatic (v0.38, sliding
359 window size 4bp, Q15, minimum read length:75bp) [62]. Metagenomic sequence reads were mapped to VIRGO
360 using bowtie (v1; parameters: -p 16 -l 25 --fullref --chunkmbs 512 --best --strata -m 20 --suppress 2,4,5,6,7,8),
361 producing a taxonomic and gene annotation for each read. Samples with fewer than 100,000 mapped reads were
362 removed from the analysis (n=59). The number of reads mapped to a gene was multiplied by the read length (150
363 bp) and divided by the gene length to produce a coverage value for each gene. Conserved domain and motif searches
364 were performed with CD-SEARCH and the Conserved Domain Database (CDD), using an e-value threshold of 10^{-4} .
365 The taxonomic composition table generated using VIRGO were run through the vaginal CST classifier VALENCIA
366 [11].

367

368 **Metagenomic Subspecies.** For each species, a presence/absence matrix was constructed from a metagenome which
369 included all genes with at least 0.5X coverage after normalizing for gene length. Metagenomic subspecies were
370 generated for species present (>75% estimated median number of genes encoded in reference genomes from the
371 Genome Taxonomy Database [63], see **Table S4 GENOME SIZES**) in >20 samples using binary gene counts and
372 hierarchical clustering with Ward linkage of sample Jaccard distances calculated using the vegdist function from the
373 vegan package (v2.5-5) [64] in R (v. 3.5.2). Clusters were defined using the dynamic hybrid tree cut method (v.1.62-
374 1) and minClusterSize = 2 [65]. Clusters were tested for associations with low species coverage using logistic
375 regression in which the mgSs was the binary outcome, the \log_{10} -transformed coverage of the species was the
376 predictor, and subject ID was used as a nested random effect which accounted for multiple samples from the same
377 subject and variations due to different source studies. Heatmaps of gene presence/absence were constructed for each
378 species using the gplots package heatmap.2 function [66] (**Supplemental File 4**).

379

380 **Metagenomic CSTs.** Using gene abundance information (normalized by gene length and sequencing depth), we
381 estimated the proportion of vaginal species in each sample. For species that were sub-divided into mgSs, the mgSs
382 proportion in a sample was equal to the proportion of the species in that sample. When a species was present in a
383 sample but with too few genes present to constitute a mgSs (<75% estimated median number of genes encoded in
384 reference genomes), it was labeled as “mgSs 0”. Samples in the resulting compositional table were hierarchically
385 clustered using Jensen-Shannon distances. Clusters were defined using the dynamic hybrid tree cut method (v.1.62-
386 1) [65]. A heatmap for metagenomic CSTs was produced using the gplots package heatmap.2 function (**Figure 1**)
387 [66]. For participants in the HMP cohort who contributed longitudinal samples, the Yue-Clayton theta was measured
388 to define microbiota stability for each subject [67]. Average per-subject stability thetas were plotted for each
389 mgCST.

390
391 **Estimating the number of *L. crispatus* strains.** The number of *L. crispatus* strains in a mgSs was estimated using a
392 pangene accumulation curve which was generated by mapping the gene contents of publicly available isolate
393 genome sequences (**Supplemental File 5**) to VIRGO (blastn, threshold: 90% identity, 70% coverage). Bootstrap
394 (n=100) combinations of N (N=1 to 61) isolates were selected and the number of unique *L. crispatus* Vaginal
395 Orthologous Groups [VOGs; provided in the VIRGO output[29] encoded in their genomes was determined. An
396 exponential curve relating the number of isolates to the number of VOGs detected was then fit to the resulting data
397 and produced the equation: $Y=2057N^{0.14}$ where Y is the number of *L. crispatus* VOGs detected, and N is the
398 estimated number of strains. This equation was then used to estimate the number of *L. crispatus* strain's detected in a
399 metagenome based on the observed number of *L. crispatus* VOGs in each metagenome.

400
401 **Statistical analysis of the association between mgCST and age, race, Nugent score, vaginal pH, and BV.** For
402 those samples with race, age category, Nugent score category, vaginal pH category, or Amsel-BV diagnoses
403 information (**TABLE**), the Cochran-Mantel-Haenszel Chi-Squared Test (CMH test, “mantelhaen.test” from the
404 samplesizeCMH R package, v 0.0.0, github.com/pegeler/samplesizeCMH) was used to determine associations with
405 mgCSTs while accounting for source study (the confounding variable). The CMH test evaluates associations
406 between two binary variables (*i.e.*, “mgCST X or not” and “high Nugent score or not”). Tests were done at the

407 subject level; if a subject had more than one sample and both samples were the same mgCST, only one sample was
408 used, but if the mgCSTs differed, the samples were included in each.

409

410 **Construction of the random forests for mgSs classification**

411 We constructed random forests for classification of mgSs using the R package randomForestSRC v2.12.1R [68]. For
412 mgSs, a random forest was built for each species (n=28) where the training data contained presence/absence values
413 of genes. Gene presence was defined as above for mgSs. We implemented random forest classification analysis with
414 all predictors included in a single model. For each mgSs random forest, predictors were all genes in a species. Ten-
415 fold cross-validation (90% of data as training, 10% as testing) was performed wherein each training set was used to
416 build and tune a random forest model using tune “tune.rfsrc”. A random forest model using optimal parameters was
417 then used to predict mgSs classifications for the test set and out-of-bag error estimates (misclassification error) are
418 reported. The overall misclassification error is the average misclassification error from each fold and the “correct”
419 assignment is based on original hierarchical clustering assignment. The final models included all data and the
420 optimal tuning parameters determined for that species.

421

422 **Construction of the a nearest centroid classifier for mgCSTs**

423 Using mgCSTs as defined above, reference centroids were produced using the mean relative abundances of each
424 mgSs in a mgCST. For classification, the similarity of a sample to the reference centroids is determined using Yue-
425 Clayton’s θ [67]. Ten-fold cross validation was applied wherein each training set was used to build “reference”
426 centroids and each test set was used for assignment. The misclassification error was determined by subtracting the
427 number of correct assignments (based on original hierarchical clustering assignment) divided by the total number of
428 assignments from 1. The overall misclassification error is the average of misclassification error from each fold.

429

430 **Running the mgCST classifier**

431 The required inputs are direct outputs from VIRGO [29] and include the taxonomic abundance table
432 (“summary.Abundance.txt”) and gene abundance table (“summary.NR.abundance.txt”). It is *imperative* that
433 taxonomic and gene column headings match those output by VIRGO. The expected output is a count table with
434 samples as rows, taxa as columns, and counts normalized by gene length as values. Additional columns indicate the

435 sample mgCST classification and the Yue-Clayton similarity score for all 26 mgCSTs. A heatmap is also produced
436 showing taxon relative abundances in samples, where samples are labeled with assigned mgCSTs Substantial
437 differences may indicate either an incongruence in taxonomic or gene names or the need for an additional mgCST.
438 The classifier is contained in an R script, which is available at <https://github.com/ravel-lab/mgCST-classifier>.
439
440 All bioinformatic and statistical analyses are available in R Markdown notebooks (**Supplemental File 7**
441 **mgCST_paper_bioinformatics.Rmd** and **Supplemental File 8 mgCST_paper_stats.Rmd**)
442

443 REFERENCES

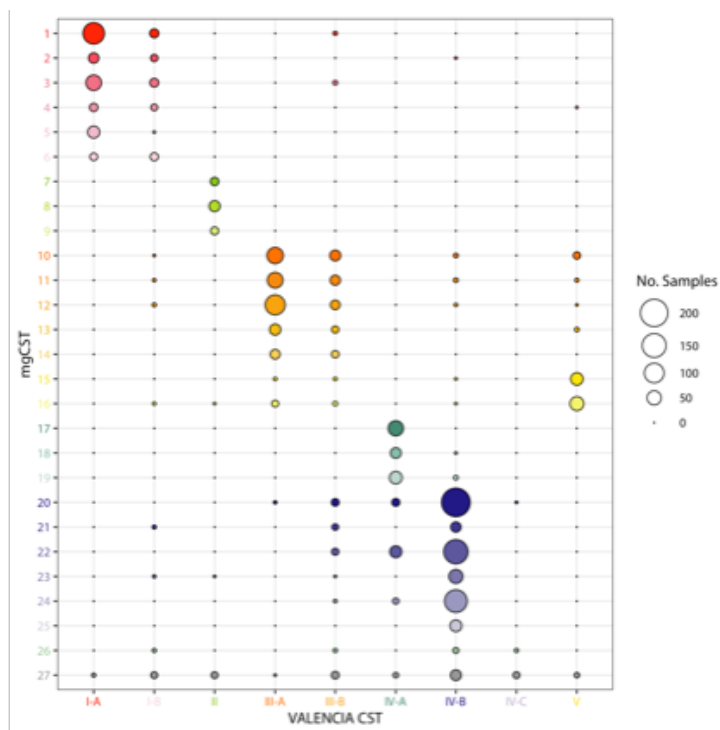
- 444 1. Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SSK, McCulle SL, et al.: Vaginal microbiome of
445 reproductive-age women. In: Proc Natl Acad Sci USA. vol. 108 Suppl 1; 2011: 4680-7.
- 446 2. O'Hanlon DE, Come RA, Moench TR. Vaginal pH measured in vivo: lactobacilli determine pH and lactic
447 acid concentration. *Bmc Microbiol.* 2019;19(1):13; doi: 10.1186/s12866-019-1388-8.
- 448 3. Gong Z, Luna Y, Yu P, Fan H. Lactobacilli inactivate *Chlamydia trachomatis* through lactic acid but not
449 H₂O₂. *PLoS One.* 2014;9(9):e107758; doi: 10.1371/journal.pone.0107758.
- 450 4. Witkin SS, Mendes-Soares H, Linhares IM, Jayaram A, Ledger WJ, Forney LJ. Influence of vaginal
451 bacteria and D- and L-lactic acid isomers on vaginal extracellular matrix metalloproteinase inducer:
452 implications for protection against upper genital tract infections. *mBio.* 2013;4(4); doi:
453 10.1128/mBio.00460-13.
- 454 5. Ravel J, Brotman RM. Translating the vaginal microbiome: gaps and challenges. *Genome Med.*
455 2016;8(1):35; doi: 10.1186/s13073-016-0291-2.
- 456 6. Amabebe E, Anumba DOC. The Vaginal Microenvironment: The Physiologic Role of Lactobacilli. *Front*
457 *Med (Lausanne).* 2018;5:181; doi: 10.3389/fmed.2018.00181.
- 458 7. Ma B, Forney LJ, Ravel J. Vaginal microbiome: rethinking health and disease. *Annu Rev Microbiol.*
459 2012;66:371-89; doi: 10.1146/annurev-micro-092611-150157.
- 460 8. Boskey ER, Cone RA, Whaley KJ, Moench TR: Origins of vaginal acidity: high D/L lactate ratio is
461 consistent with bacteria being the primary source. In: *Hum Reprod.* vol. 16: Oxford University Press; 2001:
462 1809-13.
- 463 9. Nunn KL, Wang YY, Harit D, Humphrys MS, Ma B, Cone R, et al. Enhanced Trapping of HIV-1 by
464 Human Cervicovaginal Mucus Is Associated with *Lactobacillus crispatus*-Dominant Microbiota. *mBio.*
465 2015;6(5):e01084-15; doi: 10.1128/mBio.01084-15.
- 466 10. Edwards VL, Smith SB, McComb EJ, Tamarelle J, Ma B, Humphrys MS, et al. The Cervicovaginal
467 Microbiota-Host Interaction Modulates *Chlamydia trachomatis* Infection. *mBio.* 2019;10(4); doi:
468 10.1128/mBio.01548-19.
- 469 11. France MT, Ma B, Gajer P, Brown S, Humphrys MS, Holm JB, et al. VALENCIA: a nearest centroid
470 classification method for vaginal microbial communities based on composition. *Microbiome.*
471 2020;8(1):166; doi: 10.1186/s40168-020-00934-6.
- 472 12. Brotman RM, Bradford LL, Conrad M, Gajer P, Ault K, Peralta L, et al. Association between *Trichomonas*
473 *vaginalis* and vaginal bacterial community composition among reproductive-age women. *Sexually*
474 *transmitted diseases.* 2012;39(10):807-12; doi: 10.1097/OLQ.0b013e3182631c79.
- 475 13. Mehta SD, Donovan B, Weber KM, Cohen M, Ravel J, Gajer P, et al. The vaginal microbiota over an 8- to
476 10-year period in a cohort of HIV-infected and HIV-uninfected women. *PLoS One.* 2015;10(2):e0116894;
477 doi: 10.1371/journal.pone.0116894.
- 478 14. Dunlop AL, Satten GA, Hu YJ, Knight AK, Hill CC, Wright ML, et al. Vaginal Microbiome Composition
479 in Early Pregnancy and Risk of Spontaneous Preterm and Early Term Birth Among African American
480 Women. *Front Cell Infect Microbiol.* 2021;11:641005; doi: 10.3389/fcimb.2021.641005.
- 481 15. Price JT, Vwalika B, Hobbs M, Nelson JAE, Stringer EM, Zou F, et al. Highly diverse anaerobe-
482 predominant vaginal microbiota among HIV-infected pregnant women in Zambia. *PLoS One.*
483 2019;14(10):e0223128; doi: 10.1371/journal.pone.0223128.
- 484 16. Gosmann C, Anahtar MN, Handley SA, Farcasanu M, Abu-Ali G, Bowman BA, et al. *Lactobacillus*-
485 Deficient Cervicovaginal Bacterial Communities Are Associated with Increased HIV Acquisition in Young
486 South African Women. *Immunity.* 2017;46(1):29-37; doi: 10.1016/j.immuni.2016.12.013.
- 487 17. Atashili J, Poole C, Ndumbe PM, Adimora AA, Smith JS. Bacterial vaginosis and HIV acquisition: a meta-
488 analysis of published studies. *AIDS.* 2008;22(12):1493-501; doi: 10.1097/QAD.0b013e3283021a37.
- 489 18. Elovitz MA, Gajer P, Riis V, Brown AG, Humphrys MS, Holm JB, et al. Cervicovaginal microbiota and
490 local immune response modulate the risk of spontaneous preterm delivery. *Nat Commun.* 2019;10(1):1305;
491 doi: 10.1038/s41467-019-09285-9.
- 492 19. Borgdorff H, Tsvitshivadze E, Verhelst R, Marzorati M, Jurriaans S, Ndayisaba GF, et al. *Lactobacillus*-
493 dominated cervicovaginal microbiota associated with reduced HIV/STI prevalence and genital HIV viral
494 load in African women. *ISME J.* 2014;8(9):1781-93; doi: 10.1038/ismej.2014.26.
- 495 20. Tamarelle J, Thiebaut ACM, de Barbeyrac B, Bebear C, Ravel J, Delarocque-Astagneau E. The vaginal
496 microbiota and its association with human papillomavirus, *Chlamydia trachomatis*, *Neisseria gonorrhoeae*

- 497 and *Mycoplasma genitalium* infections: a systematic review and meta-analysis. *Clin Microbiol Infect.*
498 2019;25(1):35-47; doi: 10.1016/j.cmi.2018.04.019.
- 499 21. Amsel R, Totten PA, Spiegel CA, Chen KC, Eschenbach D, Holmes KK. Nonspecific vaginitis: diagnostic
500 criteria and microbial and epidemiologic associations. *The American journal of medicine.* 1983;74(1):14-
501 22.
- 502 22. Scharbo-Dehaan M, Anderson DG. The CDC 2002 guidelines for the treatment of sexually transmitted
503 diseases: implications for women's health care. *J Midwifery Womens Health.* 2003;48(2):96-104; doi:
504 10.1016/s1526-9523(02)00416-6.
- 505 23. Bilardi JE, Walker S, Temple-Smith M, McNair R, Mooney-Somers J, Bellhouse C, et al. The burden of
506 bacterial vaginosis: women's experience of the physical, emotional, sexual and social impact of living with
507 recurrent bacterial vaginosis. *PLoS One.* 2013;8(9):e74378; doi: 10.1371/journal.pone.0074378.
- 508 24. Klebanoff MA, Schwebke JR, Zhang J, Nansel TR, Yu KF, Andrews WW. Vulvovaginal symptoms in
509 women with bacterial vaginosis. *Obstet Gynecol.* 2004;104(2):267-72; doi:
510 10.1097/01.AOG.0000134783.98382.b0.
- 511 25. Nugent RP, Krohn MA, Hillier SL. Reliability of diagnosing bacterial vaginosis is improved by a
512 standardized method of gram stain interpretation. *Journal of clinical microbiology.* 1991;29(2):297-301.
- 513 26. McKinnon LR, Achilles SL, Bradshaw CS, Burgener A, Crucitti T, Fredricks DN, et al. The Evolving
514 Facets of Bacterial Vaginosis: Implications for HIV Transmission. *AIDS Res Hum Retroviruses.*
515 2019;35(3):219-28; doi: 10.1089/AID.2018.0304.
- 516 27. Sela U, Euler CW, Correa da Rosa J, Fischetti VA. Strains of bacterial species induce a greatly varied acute
517 adaptive immune response: The contribution of the accessory genome. *PLoS Pathog.* 2018;14(1):e1006726;
518 doi: 10.1371/journal.ppat.1006726.
- 519 28. Douillard FP, Ribbera A, Kant R, Pietila TE, Jarvinen HM, Messing M, et al. Comparative genomic and
520 functional analysis of 100 *Lactobacillus rhamnosus* strains and their comparison with strain GG. *PLoS*
521 *Genet.* 2013;9(8):e1003683; doi: 10.1371/journal.pgen.1003683.
- 522 29. Ma B, France MT, Crabtree J, Holm JB, Humphrys MS, Brotman RM, et al. A comprehensive non-
523 redundant gene catalog reveals extensive within-community intraspecies diversity in the human vagina. *Nat*
524 *Commun.* 2020;11(1):940; doi: 10.1038/s41467-020-14677-3.
- 525 30. Tortelli BA, Lewis AL, Fay JC. The structure and diversity of strain-level variation in vaginal bacteria.
526 *Microb Genom.* 2021;7(3); doi: 10.1099/mgen.0.000543.
- 527 31. O'Hanlon DE, Moench TR, Cone RA. Vaginal pH and microbicidal lactic acid when lactobacilli dominate
528 the microbiota. *PLoS One.* 2013;8(11):e80074; doi: 10.1371/journal.pone.0080074.
- 529 32. Tachedjian G, Aldunate M, Bradshaw CS, Cone RA. The role of lactic acid production by probiotic
530 *Lactobacillus* species in vaginal health. *Res Microbiol.* 2017;168(9-10):782-92; doi:
531 10.1016/j.resmic.2017.04.001.
- 532 33. Kochhar S, Hottinger H, Chuard N, Taylor PG, Atkinson T, Scawen MD, et al. Cloning and overexpression
533 of *Lactobacillus helveticus* D-lactate dehydrogenase gene in *Escherichia coli*. *Eur J Biochem.*
534 1992;208(3):799-805; doi: 10.1111/j.1432-1033.1992.tb17250.x.
- 535 34. Petrova MI, Reid G, Vaneechoutte M, Lebeer S. *Lactobacillus iners*: Friend or Foe? *Trends Microbiol.*
536 2017;25(3):182-91; doi: 10.1016/j.tim.2016.11.007.
- 537 35. Whelan F, Lafita A, Gilbert J, Degut C, Griffiths SC, Jenkins HT, et al. Periscope Proteins are variable-
538 length regulators of bacterial cell surface interactions. *Proc Natl Acad Sci U S A.* 2021;118(23); doi:
539 10.1073/pnas.2101349118.
- 540 36. Gaillard JL, Berche P, Frehel C, Gouin E, Cossart P. Entry of *L. monocytogenes* into cells is mediated by
541 internalin, a repeat protein reminiscent of surface antigens from gram-positive cocci. *Cell.*
542 1991;65(7):1127-41; doi: 10.1016/0092-8674(91)90009-n.
- 543 37. Rampersaud R, Planet PJ, Randis TM, Kulkarni R, Aguilar JL, Lehrer RI, et al. Inerolysin, a cholesterol-
544 dependent cytolysin produced by *Lactobacillus iners*. *J Bacteriol.* 2011;193(5):1034-41; doi:
545 10.1128/JB.00694-10.
- 546 38. Vaneechoutte M, Guschin A, Van Simaey L, Gansemans Y, Van Nieuwerburgh F, Cools P. Emended
547 description of *Gardnerella vaginalis* and description of *Gardnerella leopoldii* sp. nov., *Gardnerella piotii* sp.
548 nov. and *Gardnerella swidsinskii* sp. nov., with delineation of 13 genomic species within the genus
549 *Gardnerella*. *Int J Syst Evol Microbiol.* 2019;69(3):679-87; doi: 10.1099/ijsem.0.003200.
- 550 39. Muzny CA, Blanchard E, Taylor CM, Aaron KJ, Talluri R, Griswold ME, et al. Identification of Key
551 Bacteria Involved in the Induction of Incident Bacterial Vaginosis: A Prospective Study. *J Infect Dis.*
552 2018;218(6):966-78; doi: 10.1093/infdis/jiy243.

- 553 40. Verstraelen H, Verhelst R, Claeys G, De Backer E, Temmerman M, Vanechoutte M. Longitudinal
554 analysis of the vaginal microflora in pregnancy suggests that *L. crispatus* promotes the stability of the
555 normal vaginal microflora and that *L. gasseri* and/or *L. iners* are more conducive to the occurrence of
556 abnormal vaginal microflora. *BMC Microbiol.* 2009;9:116; doi: 10.1186/1471-2180-9-116.
- 557 41. Ravel J, Brotman RM, Gajer P, Ma B, Nandy M, Fadrosch DW, et al.: Daily temporal dynamics of vaginal
558 microbiota before, during and after episodes of bacterial vaginosis. In: *Microbiome*. vol. 1: BioMed
559 Central; 2013: 29.
- 560 42. Ferris MJ, Norori J, Zozaya-Hinchliffe M, Martin DH: Cultivation-Independent Analysis of Changes in
561 Bacterial Vaginosis Flora Following Metronidazole Treatment. In: *Journal of Clinical Microbiology*. vol.
562 45; 2007: 1016-8.
- 563 43. Peebles K, Velloza J, Balkus JE, McClelland RS, Barnabas RV. High Global Burden and Costs of Bacterial
564 Vaginosis: A Systematic Review and Meta-Analysis. *Sex Transm Dis.* 2019;46(5):304-11; doi:
565 10.1097/OLQ.0000000000000972.
- 566 44. Macklaim JM, Fernandes AD, Di Bella JM, Hammond J-A, Reid G, Gloor GB: Comparative meta-RNA-
567 seq of the vaginal microbiota and differential expression by *Lactobacillus iners* in health and dysbiosis. In:
568 *Microbiome*. vol. 1; 2013: 12.
- 569 45. Canova MJ, Molle V. Bacterial serine/threonine protein kinases in host-pathogen interactions. *J Biol Chem.*
570 2014;289(14):9473-9; doi: 10.1074/jbc.R113.529917.
- 571 46. Tamarelle J, Ma B, Gajer P, Humphrys MS, Terplan M, Mark KS, et al. Nonoptimal Vaginal Microbiota
572 After Azithromycin Treatment for Chlamydia trachomatis Infection. *J Infect Dis.* 2020;221(4):627-35; doi:
573 10.1093/infdis/jiz499.
- 574 47. Potter RF, Burnham CD, Dantas G. In Silico Analysis of Gardnerella Genomespecies Detected in the
575 Setting of Bacterial Vaginosis. *Clin Chem.* 2019;65(11):1375-87; doi: 10.1373/clinchem.2019.305474.
- 576 48. Ksiezarek M, Ugarcina-Perovic S, Rocha J, Grosso F, Peixe L. Long-term stability of the urogenital
577 microbiota of asymptomatic European women. *Bmc Microbiol.* 2021;21(1):64; doi: 10.1186/s12866-021-
578 02123-3.
- 579 49. Turner E, Sobel JD, Akins RA. Prognosis of recurrent bacterial vaginosis based on longitudinal changes in
580 abundance of *Lactobacillus* and specific species of *Gardnerella*. *PLoS One.* 2021;16(8):e0256445; doi:
581 10.1371/journal.pone.0256445.
- 582 50. Zozaya-Hinchliffe M, Lillis R, Martin DH, Ferris MJ: Quantitative PCR Assessments of Bacterial Species
583 in Women with and without Bacterial Vaginosis. In: *Journal of Clinical Microbiology*. vol. 48; 2010: 1812-
584 9.
- 585 51. Janulaitiene M, Paliulyte V, Grinceviciene S, Zakareviciene J, Vladisauskiene A, Marcinkute A, et al.
586 Prevalence and distribution of *Gardnerella vaginalis* subgroups in women with and without bacterial
587 vaginosis. *BMC Infect Dis.* 2017;17(1):394; doi: 10.1186/s12879-017-2501-y.
- 588 52. Gelber SE, Aguilar JL, Lewis KL, Ratner AJ. Functional and phylogenetic characterization of Vaginolysin,
589 the human-specific cytolysin from *Gardnerella vaginalis*. *J Bacteriol.* 2008;190(11):3896-903; doi:
590 10.1128/JB.01965-07.
- 591 53. Pleckaityte M, Janulaitiene M, Lasickiene R, Zvirbliene A. Genetic and biochemical diversity of
592 *Gardnerella vaginalis* strains isolated from women with bacterial vaginosis. *FEMS Immunol Med*
593 *Microbiol.* 2012;65(1):69-77; doi: 10.1111/j.1574-695X.2012.00940.x.
- 594 54. Yeoman CJ, Yildirim S, Thomas SM, Durkin AS, Torralba M, Sutton G, et al. Comparative genomics of
595 *Gardnerella vaginalis* strains reveals substantial differences in metabolic and virulence potential. *PLoS*
596 *One.* 2010;5(8):e12411; doi: 10.1371/journal.pone.0012411.
- 597 55. Koumans EH, Sternberg M, Bruce C, McQuillan G, Kendrick J, Sutton M, et al. The prevalence of
598 bacterial vaginosis in the United States, 2001-2004; associations with symptoms, sexual behaviors, and
599 reproductive health. *Sex Transm Dis.* 2007;34(11):864-9; doi: 10.1097/OLQ.0b013e318074e565.
- 600 56. Workowski KA, Bachmann LH, Chan PA, Johnston CM, Muzny CA, Park I, et al. Sexually Transmitted
601 Infections Treatment Guidelines, 2021. *MMWR Recomm Rep.* 2021;70(4):1-187; doi:
602 10.15585/mmwr.rr7004a1.
- 603 57. Brotman RM, Ravel J, Cone RA, Zenilman JM. Rapid fluctuation of the vaginal microbiota measured by
604 Gram stain analysis. *Sex Transm Infect.* 2010;86(4):297-302; doi: 10.1136/sti.2009.040592.
- 605 58. Gajer P, Brotman RM, Bai G, Sakamoto J, Schütte UME, Zhong X, et al.: Temporal Dynamics of the
606 Human Vaginal Microbiota. In: *Sci Transl Med*. vol. 4: American Association for the Advancement of
607 Science; 2012: 132ra52-ra52.

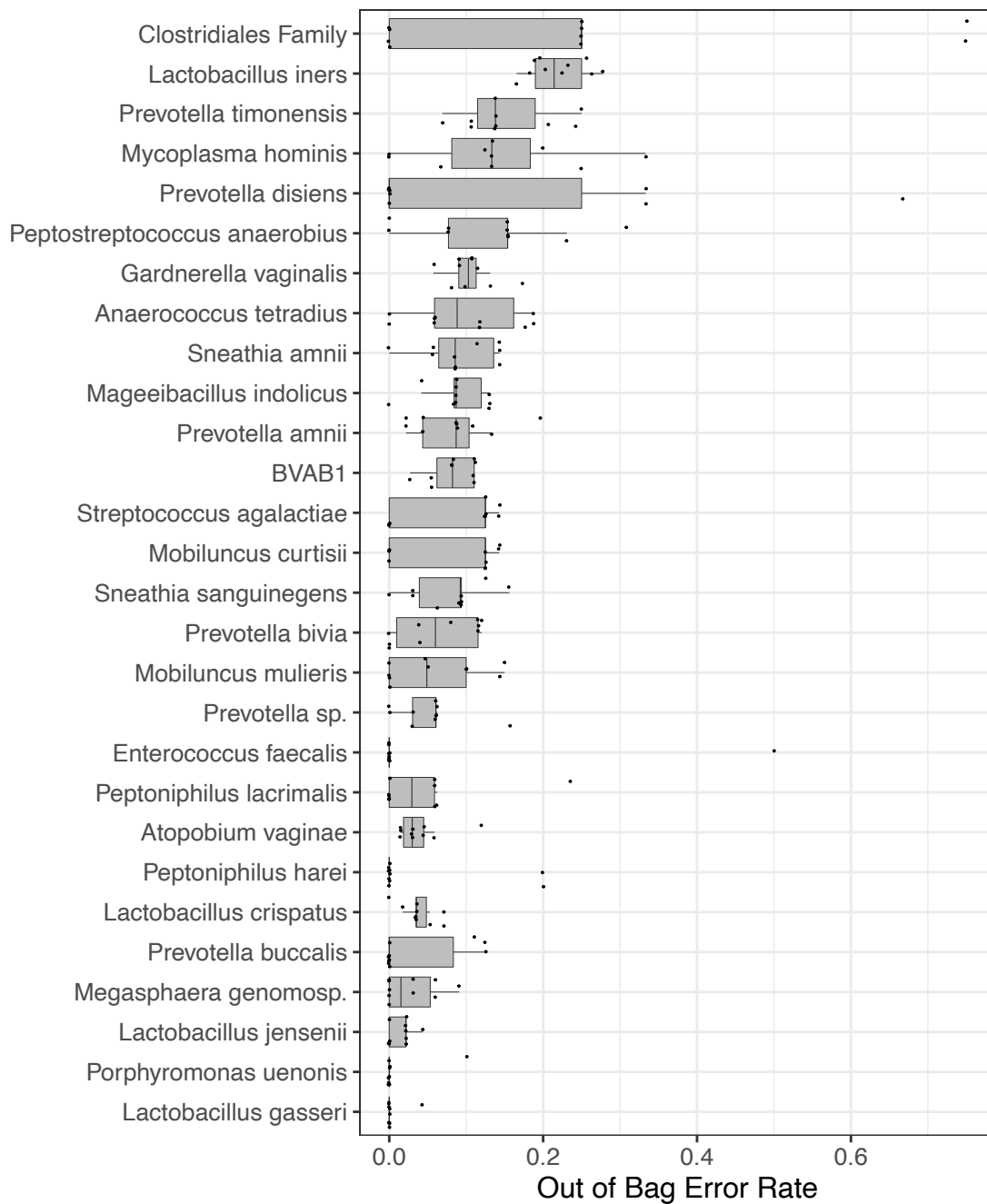
- 608 59. Munoz A, Hayward MR, Bloom SM, Rocafort M, Ngcapu S, Mafunda NA, et al. Modeling the temporal
609 dynamics of cervicovaginal microbiota identifies targets that may promote reproductive health.
610 *Microbiome*. 2021;9(1):163; doi: 10.1186/s40168-021-01096-9.
611 60. Gevers D, Knight R, Petrosino JF, Huang K, McGuire AL, Birren BW, et al. The Human Microbiome
612 Project: a community resource for the healthy human microbiome. 2012.
613 61. Li F, Chen C, Wei W, Wang Z, Dai J, Hao L, et al. The metagenome of the female upper reproductive tract.
614 *Gigascience*. 2018;7(10); doi: 10.1093/gigascience/giy107.
615 62. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data.
616 *Bioinformatics*. 2014;30(15):2114-20; doi: 10.1093/bioinformatics/btu170.
617 63. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, et al. A standardized
618 bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat Biotechnol*.
619 2018;36(10):996-1004; doi: 10.1038/nbt.4229.
620 64. Oksanen J, Kindt R, Legendre P, O'Hara B, Stevens MHH, Oksanen MJ, et al. The vegan package.
621 *Community ecology package*. 2007;10:631-7.
622 65. Langfelder P, Zhang B, Horvath S. Defining clusters from a hierarchical cluster tree: the Dynamic Tree Cut
623 package for R. *Bioinformatics*. 2008;24(5):719-20; doi: 10.1093/bioinformatics/btm563.
624 66. Warnes MGR, Bolker B, Bonebakker L, Gentleman R. Package 'gplots'. *Various R Programming Tools for*
625 *Plotting Data*. 2016.
626 67. Yue JC, Clayton MK. A similarity measure based on species proportions. *Communications in Statistics-*
627 *theory and Methods*. 2005;34(11):2123-31.
628 68. U.B. IHaK. Fast Unified Random Forests for Survival, Regression, and Classification (RF-SRC), R
629 package. 2021.
630

631 SUPPLEMENTAL FIGURE LEGENDS



633 Supplemental Figure 1. Metagenomic CSTs correspond to marker gene-based CSTs primarily through predominant
634 taxon. Dominance by mgSs is not captured through marker-based CSTs.

635

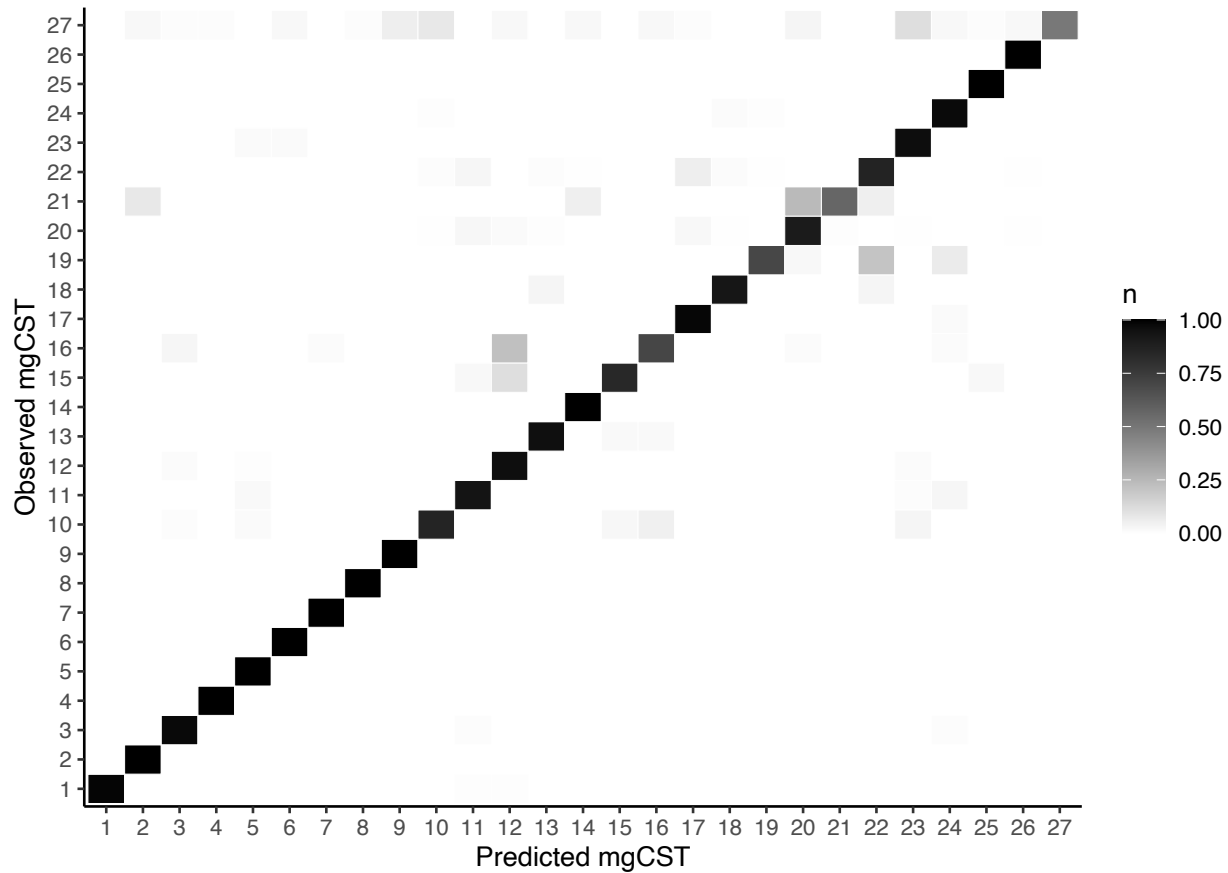


636

637 Supplemental Figure 2. Random forest misclassification error estimates from 10-fold cross-validation for each of the

638 species that contained metagenomic subspecies.

639



640

641 Supplemental Figure 3. Confusion matrix and classification error estimates from 10-fold cross-validation of a
642 nearest-centroid classifier for mgCSTs.

643

644 Supplemental File 4. Gene presence and absence heatmaps for all metagenomic subspecies.

645

646 Supplemental File 5. Gene contents of publicly available isolate metagenomes.

647

648 Supplemental Table 6. Metadata and source location for all metagenomes in this study. For Amsel-BV:

649 clinBV.asymp=Clinical diagnosis of Amsel-BV, no symptoms reported by patient; noBV=No diagnosis of Amsel-

650 BV after evaluation; clinBV.sym=Clinical diagnosis of Amsel-BV, symptoms reported by patient; NA=No clinical

651 evaluation

652

653 Supplemental_File_7_mgCST_paper_bioinformatics.Rmd. Rmarkdown notebook with code used to build
654 metagenomic subspecies and metagenomic community state types.
655
656 Supplemental_File_8_mgCST_paper_stats.Rmd. Rmarkdown notebook with code for performing all analyses and
657 generating all figures in this manuscript.
658

Dominant Species

Lactobacillus crispatus

Lactobacillus iners

"Ca. Lachnocurva vaginae"

Gardnerella

Bifidobacterium breve

Other

mgCST

1-6

7-9

10-14

15-16

17-19

20-25

26-27

Shannon Diversity

Amsel-BV

NA BV Asym BVz Sym BVz

Nugent Score

NA 0-3 4-6 7-10

Vaginal pH

NA 4 7

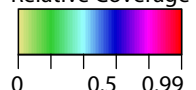
Race

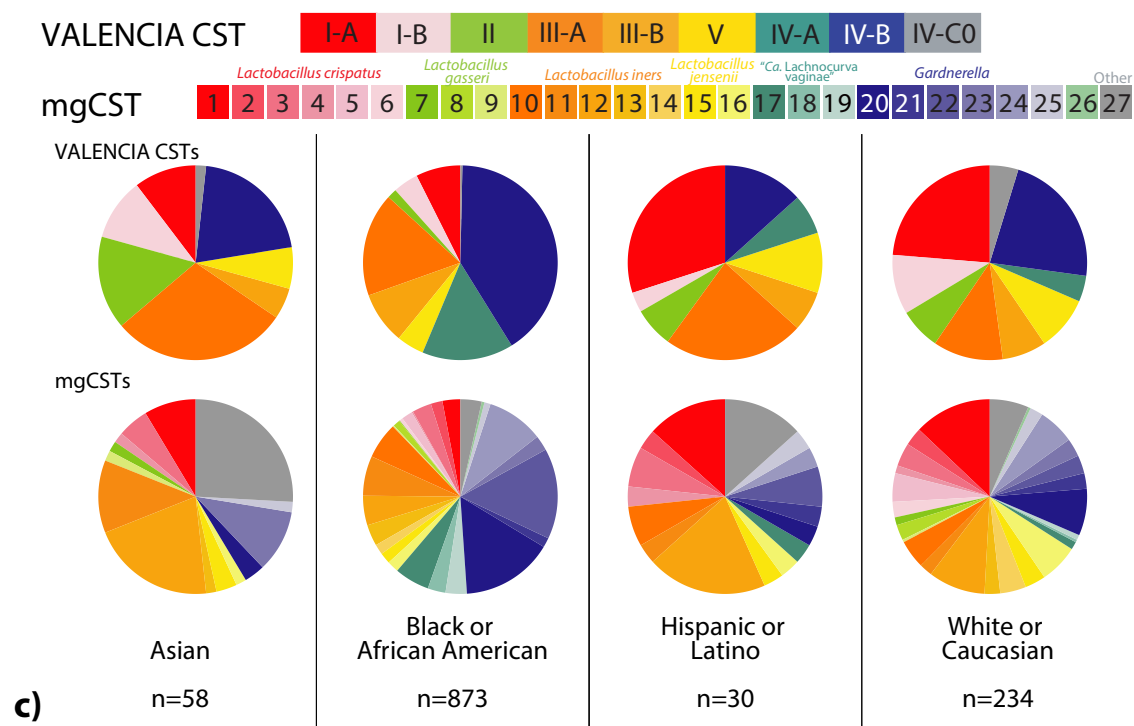
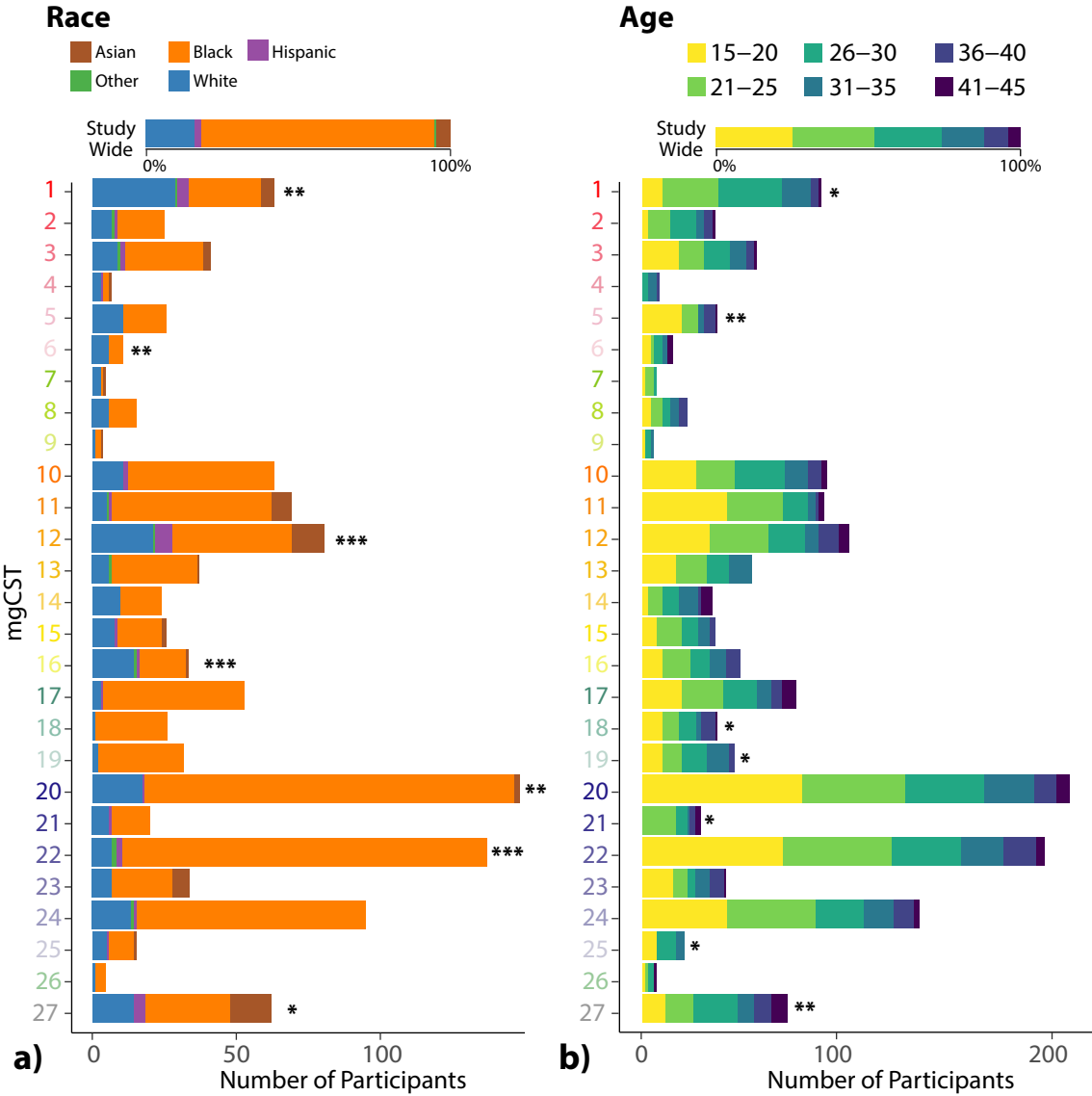
American Indian / Asian / Multi-racial / Black / Hispanic / Other / White

- Lactobacillus iners* 0
- Lactobacillus iners* 1
- Lactobacillus iners* 2
- Lactobacillus iners* 3
- Lactobacillus iners* 4
- Lactobacillus iners* 5
- Lactobacillus iners* 6
- Gardnerella vaginalis* 0
- Gardnerella vaginalis* 1
- Gardnerella vaginalis* 2
- Gardnerella vaginalis* 3
- Gardnerella vaginalis* 4
- Gardnerella vaginalis* 5
- Lactobacillus crispatus* 0
- Lactobacillus crispatus* 1
- Lactobacillus crispatus* 2
- Lactobacillus crispatus* 3
- Lactobacillus crispatus* 4
- Lactobacillus crispatus* 5
- Lactobacillus crispatus* 6
- Lactobacillus jensenii* 0
- Lactobacillus jensenii* 1
- Lactobacillus jensenii* 2
- Lactobacillus jensenii* 3
- "Ca. Lachnocurva vaginae"* 0
- "Ca. Lachnocurva vaginae"* 1
- "Ca. Lachnocurva vaginae"* 2
- "Ca. Lachnocurva vaginae"* 3
- "Ca. Lachnocurva vaginae"* 4
- "Ca. Lachnocurva vaginae"* 5
- Atopobium vaginae* 0
- Atopobium vaginae* 1
- Atopobium vaginae* 2
- Lactobacillus gasseri* 0
- Lactobacillus gasseri* 1
- Lactobacillus gasseri* 2
- Lactobacillus gasseri* 3
- Bifidobacterium breve*
- Prevotella buccalis* 0
- Prevotella buccalis* 1
- Prevotella buccalis* 2
- Prevotella buccalis* 3
- Streptococcus anginosus*
- Megasphaera genomosp.* 0
- Megasphaera genomosp.* 1
- Megasphaera genomosp.* 2
- Megasphaera genomosp.* 3
- Prevotella amnii* 0
- Prevotella amnii* 1
- Prevotella amnii* 10
- Prevotella amnii* 11
- Prevotella amnii* 2
- Prevotella amnii* 3
- Prevotella amnii* 4
- Prevotella amnii* 5
- Prevotella amnii* 6
- Prevotella amnii* 7

(Showing 54 most abundant species and subspecies out of 407)

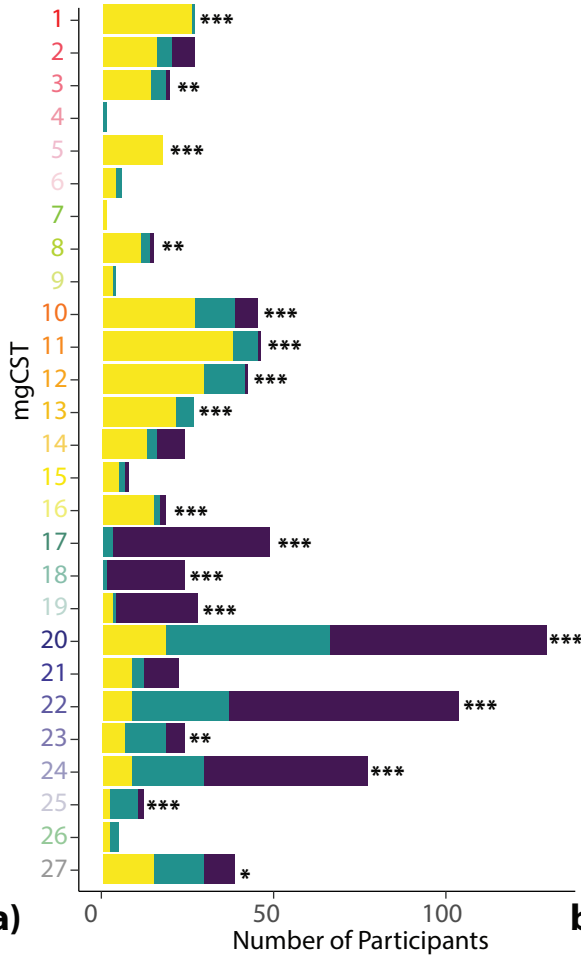
Relative Coverage





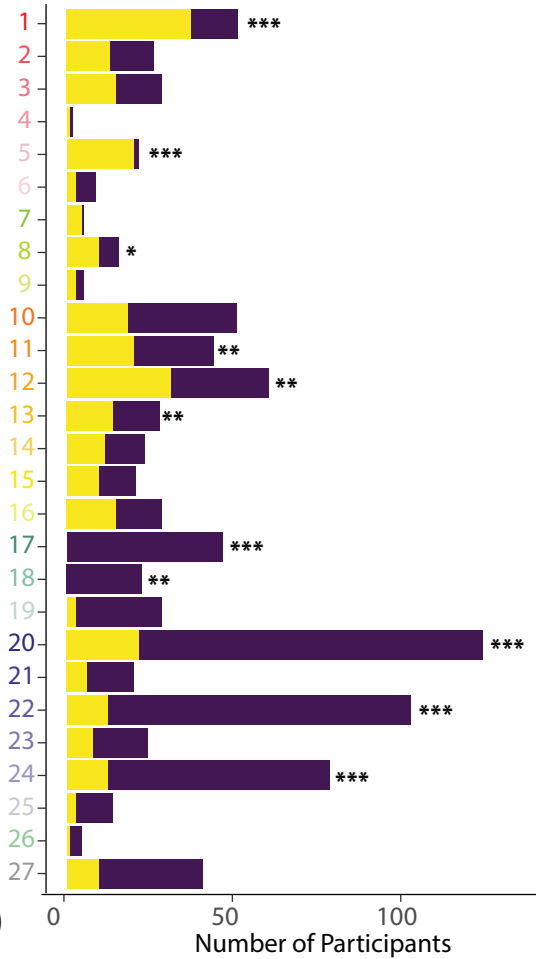
Nugent Score

Low (0–3) Intermediate (4–6) High (7–10)



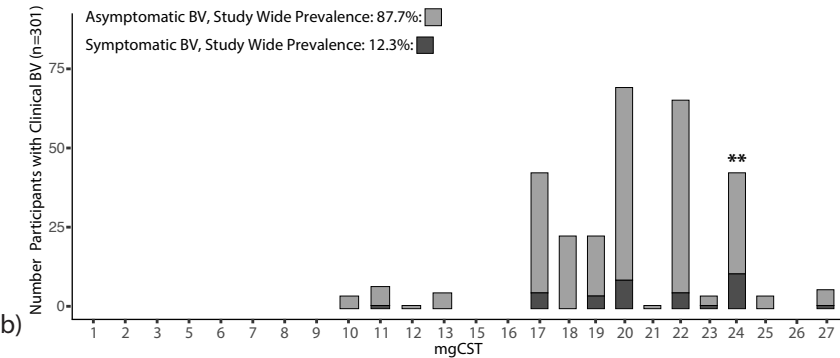
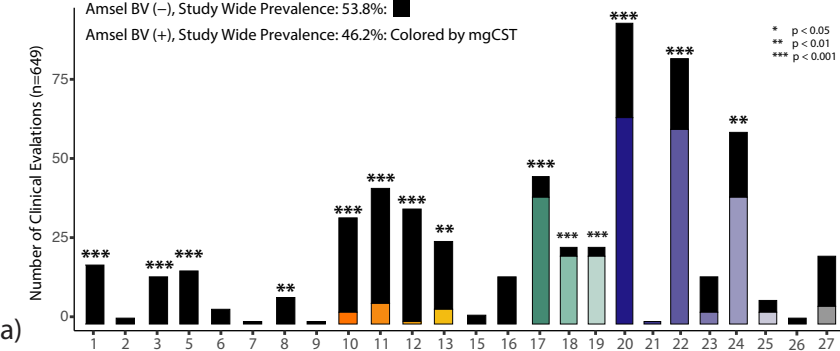
Vaginal pH

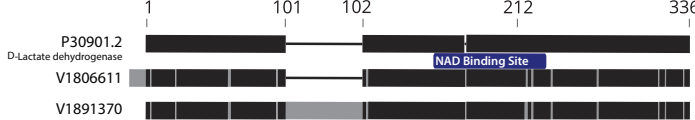
Low (< 4.5) High (≥ 4.5)



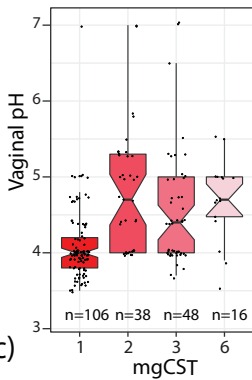
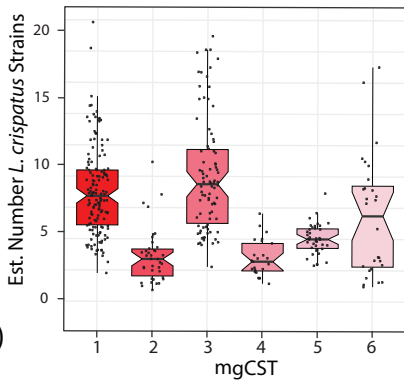
a)

b)

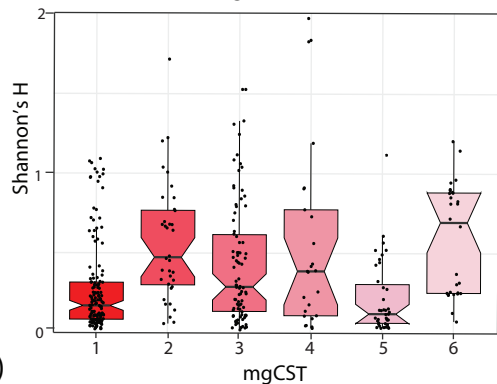




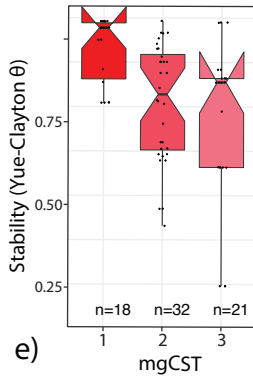
a)



b)

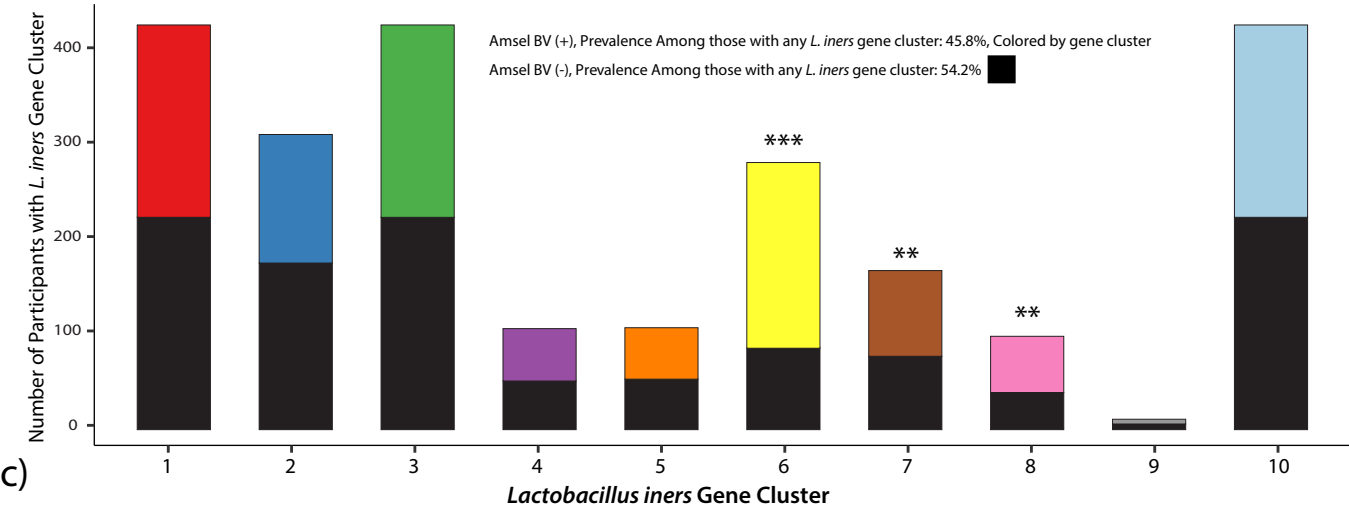
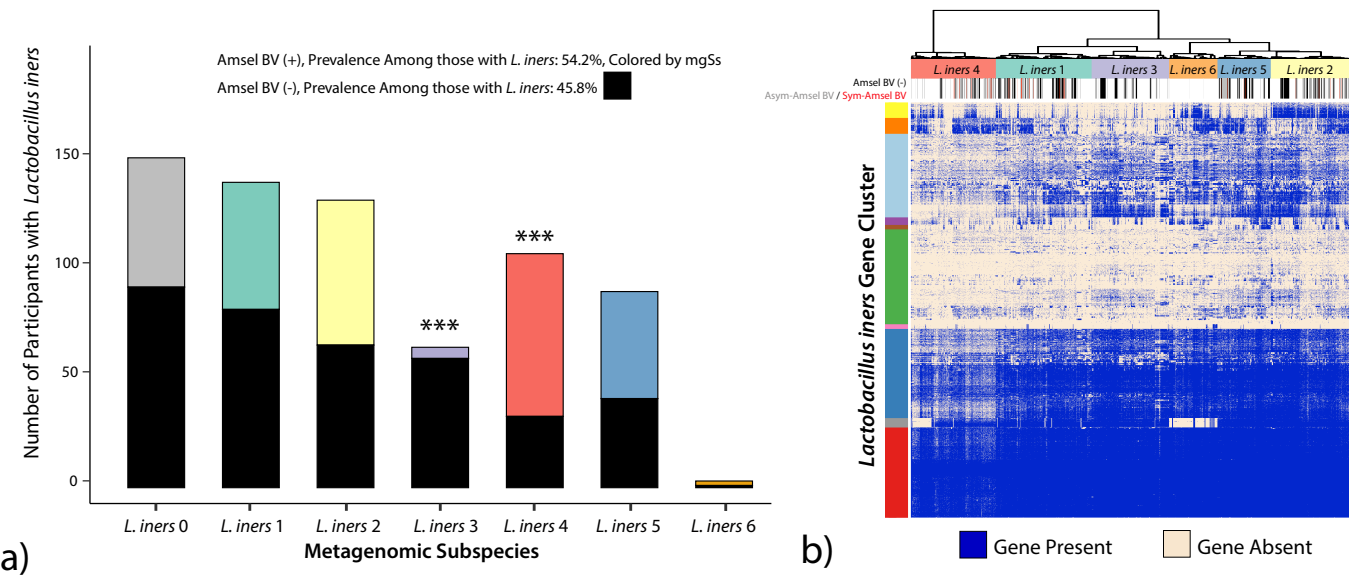


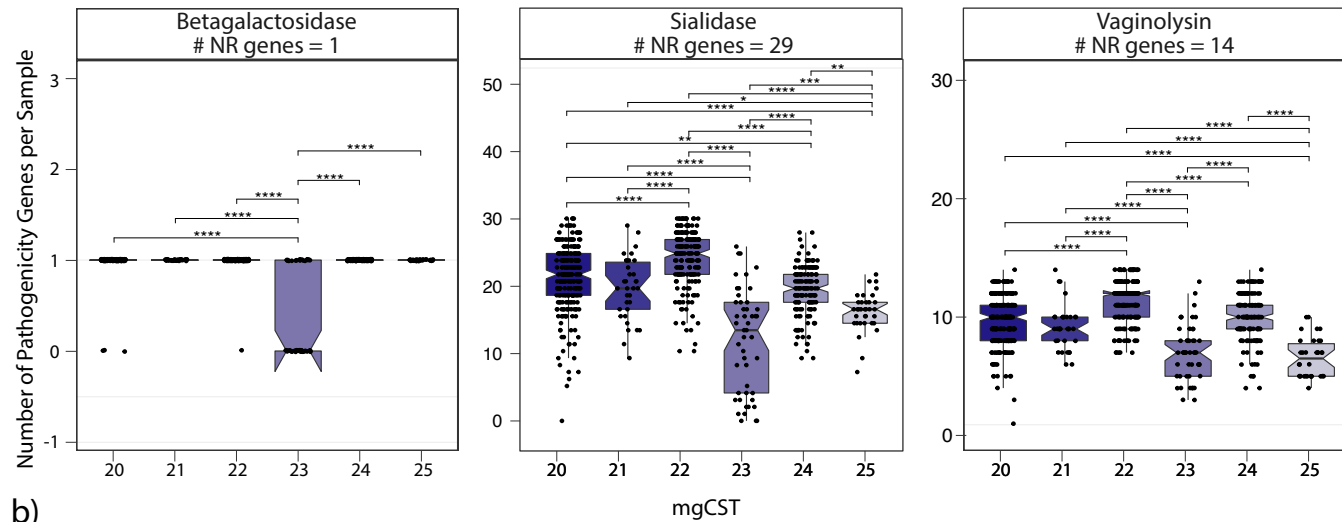
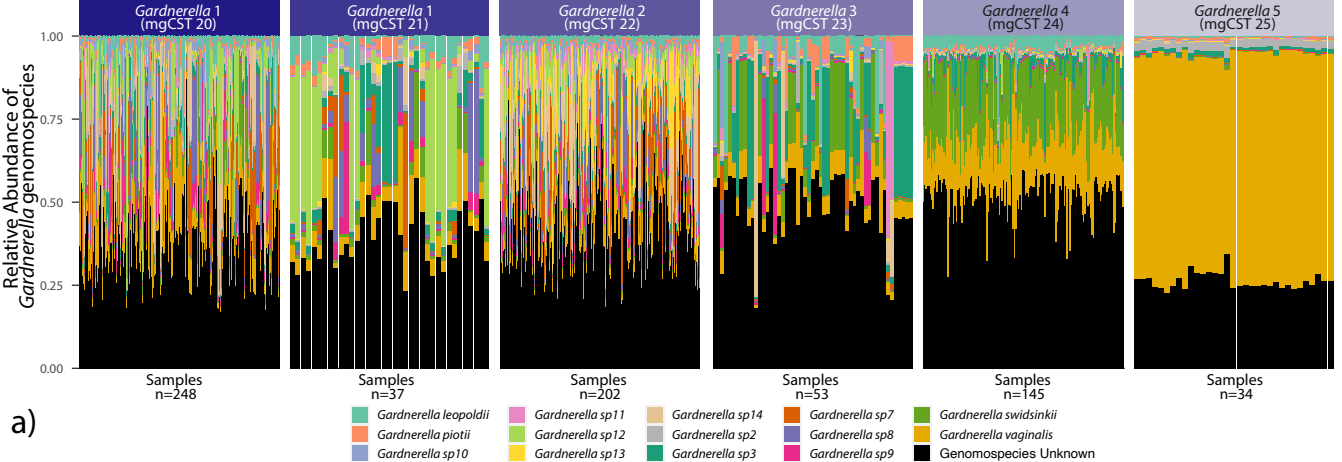
c)



d)

e)





	Number of Women	Percentage of Women	Number of Samples	Percentage of Samples
Metagenomic Data Source	1,817		1,898	
UMB-HMP	124	12.2	315	27.2
Li et al.	44	35.5	44	8.3
LSVF	585	1329.5	653	1484.1
NIH-HMP	76	13.0	174	26.6
VOMEC	40	32.6	162	93.1
VIRGO	148	370.0	342	211.1
Age Category	897		1,613	
15-20	283	31.5	410	25.3
21-25	229	25.5	436	26.9
26-30	188	21.0	362	22.3
31-35	102	11.4	223	13.7
36-40	63	7.2	123	7.7
41-45	38	3.3	67	4.1
Race	858		1,441	
Asian	54	6.3	66	4.6
Black or African American	609	71.1	968	67.2
Hispanic or Latino	39	2.2	47	3.3
Other	6	0.7	9	0.6
White or Caucasian	159	18.7	351	24.4
Age at Category	948		1,613	
0-3	469	48.5	931	57.4
4-6	194	20.0	255	15.7
7-10	305	31.5	427	26.9
Vaginal pH Category	874		1,342	
Low (pH < 4.5)	273	31.2	491	36.6
High (pH ≥ 4.5)	601	68.8	871	64.0
Anneal-EV Diagnosis	627		673	
Positive	289	46.1	298	45.8
Negative	338	53.9	365	54.2
Symptomatic Anneal-EV	289		308	
Asymptomatic	253	87.5	271	88.0
Symptomatic	36	12.5	37	12.0

Vaginal Microbiome Gene Database (VIRGO), virgo.gen.umaryland.edu [29], the University of Maryland Baltimore Human Microbiome Project (UMH-HMP), [DOI:10.1093/bioinformatics/btu004](https://doi.org/10.1093/bioinformatics/btu004), [DOI:10.1093/bioinformatics/btu005](https://doi.org/10.1093/bioinformatics/btu005), the National Institutes of Health Human Microbiome Project (NIH-HMP), [DOI:10.1093/bioinformatics/btu006](https://doi.org/10.1093/bioinformatics/btu006), Li et al. [60] [DOI:10.1093/bioinformatics/btu007](https://doi.org/10.1093/bioinformatics/btu007), the Longitudinal Study of Vaginal Flora and Incident STI (LSVF, doCaP project) [DOI:10.1093/bioinformatics/btu008](https://doi.org/10.1093/bioinformatics/btu008).

MgCST	Most Common mgSs	Most Abundant mgSs	Number of Samples	Number of Women	Median Shannon Index	Number of Samples from Metagenomic Data Source					
						UMB-HMP	Li et al.	LSVF	HMP	VIRGO	VMRC
1	<i>Lactobacillus crispatus 1</i>	<i>Lactobacillus crispatus 1</i>	143	79	0.17	20	2	21	63	15	22
2	<i>Lactobacillus crispatus 2</i>	<i>Lactobacillus crispatus 2</i>	39	26	0.47	39	0	0	0	0	0
3	<i>Lactobacillus crispatus 3</i>	<i>Lactobacillus crispatus 3</i>	83	51	0.28	9	2	15	14	22	21
4	<i>Lactobacillus crispatus 4</i>	<i>Lactobacillus crispatus 4</i>	27	12	0.39	1	1	0	0	3	22
5	<i>Lactobacillus crispatus 5</i>	<i>Lactobacillus crispatus 5</i>	37	27	0.11	1	0	19	16	1	0
6	<i>Lactobacillus crispatus 6</i>	<i>Lactobacillus crispatus 6</i>	28	13	0.69	12	0	5	2	9	0
7	<i>Lactobacillus gasseri 1</i>	<i>Lactobacillus gasseri 1</i>	16	8	0.5	0	0	1	8	1	6
8	<i>Lactobacillus gasseri 2</i>	<i>Lactobacillus gasseri 2</i>	29	17	0.73	15	0	8	0	1	5
9	<i>Lactobacillus gasseri 3</i>	<i>Lactobacillus gasseri 3</i>	14	5	0.89	6	0	0	2	0	6
10	<i>Lactobacillus iners 1</i>	<i>Lactobacillus iners 1</i>	113	76	0.7	24	0	40	2	10	37
11	<i>Lactobacillus iners 2</i>	<i>Lactobacillus iners 2</i>	95	79	0.53	11	7	47	0	28	2
12	<i>Lactobacillus iners 3</i>	<i>Lactobacillus iners 3</i>	131	92	0.44	9	10	45	19	42	6
13	<i>Lactobacillus iners 5</i>	<i>Lactobacillus iners 5</i>	45	41	0.57	1	0	29	1	13	1
14	<i>Lactobacillus iners 6</i>	<i>Lactobacillus iners 6</i>	34	25	0.8	34	0	0	0	0	0
15	<i>Lactobacillus jensenii 1</i>	<i>Lactobacillus jensenii 1</i>	44	28	0.77	8	1	3	10	18	4
16	<i>Lactobacillus jensenii 2</i>	<i>Lactobacillus jensenii 2</i>	67	39	0.71	8	0	15	15	13	16
17	"Ca." <i>Lachnocurva vaginae 1</i>	"Ca." <i>Lachnocurva vaginae 1</i>	58	57	1.48	3	0	51	1	3	0
18	"Ca." <i>Lachnocurva vaginae 1</i>	"Ca." <i>Lachnocurva vaginae 1</i>	28	27	1.57	0	0	27	0	1	0
19	"Ca." <i>Lachnocurva vaginae 1</i>	"Ca." <i>Lachnocurva vaginae 1</i>	43	36	1.91	7	0	27	0	9	0
20	<i>Gardnerella vaginalis 1</i>	<i>Gardnerella vaginalis 1</i>	250	171	1.62	90	2	98	2	38	20
21	<i>Gardnerella vaginalis 1</i>	<i>Gardnerella vaginalis 1</i>	37	21	1.97	37	0	0	0	0	0
22	<i>Gardnerella vaginalis 2</i>	<i>Prevotella amnii 4</i>	202	159	1.79	30	0	91	3	67	11
23	<i>Gardnerella vaginalis 3</i>	<i>Gardnerella vaginalis 3</i>	53	42	0.88	18	5	15	2	5	8
24	<i>Gardnerella vaginalis 4</i>	<i>Gardnerella vaginalis 4</i>	145	106	1.21	44	0	64	6	24	7
25	<i>Gardnerella vaginalis 5</i>	<i>Gardnerella vaginalis 5</i>	34	17	0.83	11	1	9	6	2	5
26	<i>Bifidobacterium breve</i>	<i>Bifidobacterium breve</i>	16	11	0.9	5	0	1	0	8	2
27	<i>Bifidobacterium dentium</i>	<i>Enterococcus faecalis 3</i>	87	76	1.78	23	13	22	2	9	18

UMB-HMP: University of Maryland Baltimore - Human Microbiome Project; Li et al. [62]; PRJEB24147; LSVF: Longitudinal Study of the Vaginal Flora; HMP: Human Microbiome Project; VIRGO: virgo.igs.umaryland.edu; VMRC: Vaginal Microbiome Research Consortium.