1 **Unique roles of vaginal *Megasphaera* phylotypes in reproductive health**

2

3 Abigail L. Glascock[1], Nicole R. Jimenez[2], Sam Boundy[2], Vishal N. Koparde[1], J. Paul

4 Brooks[3], David J. Edwards[4], Jerome F. Strauss III[5], Kimberly K. Jefferson[2], Myrna G.

5 Serrano[2], Gregory A. Buck[2,6], Vaginal Microbiome Consortium and Jennifer M.

6 Fettweis[2,5]*

7 [1]Life Sciences, Virginia Commonwealth University, Richmond, Virginia, USA

8 [2]Department of Microbiology & Immunology, Virginia Commonwealth University,

9 Richmond, Virginia, USA

10 [3]Department of Supply Chain Management and Analytics, Virginia Commonwealth

11 University, Richmond, Virginia, USA

12 [4]Department of Statistical Sciences and Operations Research, Virginia Commonwealth

13 University, Richmond, Virginia, USA

14 [5]Department of Obstetrics and Gynecology, Virginia Commonwealth University,

15 Richmond, Virginia, US

16 [6]Department of Computer Science, Virginia Commonwealth University, Richmond,

17 Virginia, USA

18 **Running Title: Vaginal *Megasphaera***

19 *Address correspondence to Jennifer M. Fettweis, jennifer.fettweis@vcuhealth.org.

20

## ABSTRACT

The composition of the human vaginal microbiome has been extensively studied and is known to influence reproductive health. However, the functional roles of individual taxa and their contributions to negative health outcomes have yet to be well characterized. Here, we examine two vaginal bacterial taxa grouped within the genus *Megasphaera* that have been previously associated with bacterial vaginosis (BV) and pregnancy complications. Phylogenetic analyses support the classification of these taxa as two distinct species. These two phylotypes, *Megasphaera* phylotype 1 (MP1) and *Megasphaera* phylotype 2 (MP2), differ in genomic structure and metabolic potential, suggestive of differential roles within the vaginal environment. Further, these vaginal taxa show evidence of genome reduction and changes in DNA base composition, which may be common features of host dependence and/or adaptation to the vaginal environment. In a cohort of 3,870 women, we observed that MP1 has a stronger positive association with bacterial vaginosis whereas MP2 was positively associated with trichomoniasis. MP1, in contrast to MP2 and other common BV-associated organisms, was not significantly excluded in pregnancy. In a cohort of 52 pregnant women, MP1 was both present and transcriptionally active in 75.4% of vaginal samples. Conversely, MP2 was largely absent in the pregnant cohort. This study provides insight into the evolutionary history, genomic potential and predicted functional role of two clinically relevant vaginal microbial taxa.

41      The vaginal microbiome is an important determinant of women's reproductive

42    health, pregnancy outcomes and neonatal health[1–6]. Optimal vaginal microbial health is

43    typically characterized by dominance of one or more lactic-acid producing species of the

44    genus *Lactobacillus* that function to lower the pH and prohibit the growth of other

45    organisms[7]. A vaginal microbiome depleted of protective vaginal lactobacilli and enriched

46    in diverse anaerobic species is often clinically diagnosed as bacterial vaginosis (BV). BV

47    is the most common vaginal condition worldwide, affecting an estimated 27% of women

48    in North America[8]. This condition has been associated with an increased risk of acquiring

49    sexually transmitted infections (STIs) as well as pregnancy complications including

50    spontaneous preterm birth[9–12]. While associations of vaginal microbial taxa with

51    reproductive health conditions such as BV are well established, the pathophysiological

52    significance of these taxa remains largely unknown. Developing a more comprehensive

53    understanding of how individual taxa contribute to negative health outcomes is essential

54    for understanding the underlying biological mechanisms and for the development of

55    effective therapeutics.

56      Here, we focus on two vaginal anaerobic taxa, *Megasphaera* phylotype 1 (MP1)

57    and *Megasphaera* phylotype 2 (MP2) and their roles in reproductive health and disease.

58    Both MP1 and MP2 have been previously associated with bacterial vaginosis across

59    multiple cohorts[3,13–15]. Due to its high specificity for the condition, MP1 has been used in

60    combination with other taxa for molecular diagnosis of BV[14,16]. MP2 was described by

61    Martin *et al.* to be more prevalent in samples collected from women with trichomoniasis,

62    suggesting the potential for divergent roles of these vaginal *Megasphaera* in disease

63    states[17]. *Megasphaera* species have also been linked to an increased risk for HIV

64    acquisition[18,19]. Given that MP1 and MP2 have been observed in the urogenital tracts of

65    adolescent males and heterosexual couples, it seems likely that these bacteria can be

66    sexually transmitted[20,21].

67    Vaginal carriage of *Megasphaera* is strongly associated with BV, and pregnant

68    women with BV have an elevated risk for spontaneous preterm birth[22]. The outcomes

69    across antibiotic intervention studies for prevention of preterm birth have been

70    inconsistent, which may be attributed in part to the significant heterogeneity in study

71    design and the choice and timing of therapeutic intervention[22]. It is now clear that there

72    are different subtypes of BV that can be stratified using molecular approaches, and some

73    subtypes of BV may be more tightly linked to preterm birth than others. Even though BV

74    has long been linked to elevated risk for preterm birth, more recent vaginal microbiome

75    studies have identified higher MP1 carriage in women who go on to deliver preterm[12,23–

76    25]. Interestingly, Mitchell *et al.* observed MP1 in samples collected from the upper genital

77    tract of women undergoing hysterectomy[26], suggesting that MP1 may be capable of

78    ascending from the vaginal environment into the upper genital tract. Together, these

79    observations suggest that MP1 can colonize the vaginal environment, ascend into the

80    upper genital tract and potentially contribute to PPROM and/or spontaneous preterm

81    birth.

82    In the current study, we use several approaches to delineate the roles of MP1 and

83    MP2 in reproductive health. These include phylogenetic analyses that probe the

84    evolutionary history of these organisms, genomic characterization that permits

85    assessment of their metabolic potential, and a study to define their individual associations

86    with demographic and clinical measures.

87    **RESULTS**

88    **Evolutionary history and genomic divergence of MP1 and MP2**

89    Genomes of three representative MP1 isolates and three representative MP2

90    isolates were analyzed to gain insight into the mechanisms underlying their colonization

91    of the human vaginal environment (Supplementary Table 1). A phylogenetic analysis of

92    145 orthologous genes using 110 genomes classified to the class Negativicutes revealed

93    that MP1 and MP2 are evolutionarily distinct and separated from the nearest

94    *Megasphaera/Anaeroglobus* clade (Fig. 1). Similar results were observed in a

95    phylogenetic analysis using 16S ribosmal RNA (rRNA) genes and the inferred topology

96    was largely reflective of niche adaptation (Fig. 2). In one case, niche-specific separation

97    did not occur; *Megasphaera* sp. BV3C16-1, which was isolated from the human vagina,

98    was grouped with oral taxa. This taxon has been reported in vaginal microbiome[24], but it

99    has been observed at low abundance and prevalence. For example, the taxon was

100   identified in five of 3,870 vaginal samples (0.12%) at a threshold of 0.01% in a cohort of

101   women enrolled through the Vaginal Human Microbiome Project (VaHMP)[27].

102   Given the significant divergence of these two vaginal phylotypes from other closely

103   related taxa, we performed a percentage of conserved proteins (POCP) analysis, a metric

104   for delineating genus boundaries[28]. The suggested cutoff for delineation of genera is a

105   POCP value of less than 50% to the genus type strain. The POCP values for members

106   of the MP1 and MP2 clade in comparison to the type strain (*Megasphaera elsdenii* DSM

107   20460) range from 49.6-52.6% (Supplementary Fig. 1, Supplementary Table 2).

108   *Anaeroglobus geminatus,* currently classified as a separate genus, had a POCP value of

109   52.5% compared to the *Megasphaera* type strain[29]. A recent study by Campbell *et al.*

110    identified Conserved Signature Indels (CSIs) and Conserved Signature Proteins (CSPs)

111    used to classify organisms to families within the class Negativicutes[30]. We identified all

112    CSIs and CSPs indicative of Veillonellaceae family genomes in MP1 and MP2 genomes

113    (Supplementary Table 3), supporting their previous placement within the Veillonellaceae

114    family. However, three of nine CSP markers specific for the class Negativicutes were

115    absent from all MP1 and MP2 genomes, indicative of genome reduction that is not

116    observed in other host-related *Megasphaera*. While biochemical analyses have yet to be

117    performed, the phylogeny, POCP analysis, loss of CSP markers, and specificity of the

118    clade to the vaginal environment could support placement of these phylotypes into a novel

119    genus of bacteria.

120        We compared the genomes of MP1 and MP2 with genomes of seven

121    *Megasphaera* isolates from human and mammalian GI tracts, the single human oral

122    *Anaeroglobus* isolate and the vaginal *Megasphaera* sp. BV3C16-1 isolate[29,31–34]. All of

123    the MP1 and MP2 isolates exhibit evidence of genome reduction with an average genome

124    size of 1.71 megabases (Mb) relative to an average genome size of 2.35 Mb for the other

125    studied *Megasphaera* and *Anaeroglobus* genomes (q=0.001, 95% CI [-0.97,-0.32],

126    Kruskal-Wallis test for differences in genome size with FDR correction). The MP1 and

127    MP2 genomes contain a predicted 1,571 protein-coding genes on average, which is

128    significantly fewer than the number of protein-coding genes for the other studied

129    *Megasphaera* and *Anaeroglobus* genomes, which contained an average of 2,116 genes

130    (q=0.00015, 95% CI [-749,-341], Kruskal-Wallis test for differences in predicted gene

131    count with FDR correction ). MP1 and MP2 also exhibit lower average GC composition

132    with an average of 42.6% compared to an average of 51.1% in the other host-associated

6

133    genomes in the *Megasphaera/Anaeroglobus* clade (q=0.0005, 95% CI [-12.17,-4.75],

134    Kruskal-Wallis test for differences in average GC composition with FDR correction) (Fig.

135    3a). Reduction in genome size and lower GC percentage has been observed in vaginal

136    strains of other bacterial taxa including *Lactobacillus and Gardnerella*, suggesting

137    reductive evolution may be a common feature of adaptation to the vaginal

138    environment[35,36].

**Taxonomic placement of MP1 and MP2 as two discrete species**

140    Similarity of the 16S rRNA gene at an identity threshold of 97% is often used to

141    delineate species. The 16S rRNA similarity between the two phylotypes is 96.3%. This

142    figure along with reports by Srinivasan *et al.*, implies that the two phylotypes are best

143    classified as distinct species based on 16S rRNA gene sequence similarity

144    (Supplementary Table 4)[37,38]. The average nucleotide identity (ANI) between MP1 and

145    MP2, which takes into account the entire nucleotide content of genomes, is 73%. This

146    figure is markedly less than the 95-96% threshold suggested for species demarcation

147    using this method (Supplementary Table 5)[39]. Our phylogenetic analyses (Fig. 1, Fig. 2)

148    reflected these findings, with MP1 and MP2 identified as sister taxa, distinct from other

149    *Megasphaera* and *Anaeroglobus* and separated by significant branch lengths, signifying

150    extensive divergence.

151    Further comparative analyses revealed that genomic synteny is conserved within

152    phylotype, with variations attributable to the presence of temperate bacteriophage.

153    However, extensive genome rearrangement was observed between MP1 and MP2

154    genomes (Fig. 3c, Supplementary Fig. 2). While a significant difference in genome size

155    was not observed between MP1 (average of 1.72 Mb) and MP2 (average of 1.70 Mb)

7

156    isolates (q=0.7497, 95% CI [-0.118, 0.152], Kruskal-Wallis test for differences in genome

157    size with FDR correction), there was an observed difference in GC composition between

158    the MP1 (average of 46.3%) and MP2 isolates (average of 39.0%) (q=0.000002, 95% CI

159    [6.95,7.61], Kruskal-Wallis test for differences in GC composition with FDR correction).

160    The two phylotypes also exhibit GC-divergent codon preference at the third position

161    (average GC composition at third position: MP1- 47%, MP2- 31%), signaling evolutionary

162    pressure for a reduction in GC composition in MP2 (Fig. 3b). The observed sequence

163    divergence, differential GC composition and codon preference, and lack of synteny

164    between MP1 and MP2 genomes provide support for the designation of the two

165    phylotypes as distinct species.

166    **Genomic evidence for niche specialization to the vaginal environment**

167         To assess differences in the predicted metabolic potential, we annotated and

168    performed metabolic reconstructions of 15 genomes including representatives of MP1,

169    MP2 and related bacterial strains classified to the *Megasphaera* and *Anaeroglobus*

170    genera. As expected, given the observed genome reduction of MP1 and MP2, many

171    metabolic pathways present among all other related taxa are absent in the MP1/MP2

172    clade (Supplementary Table 6). MP1 and MP2 are predicted to lack genes conserved in

173    other *Megasphaera* and *Anaeroglobus* genomes that function to transport putrescine and

174    spermidine, metabolize nitrogen, produce selenocysteine and transport and modify the

175    metals nickel and molybdenum. Thus, these organisms may have evolved to rely on

176    synergy with the host and/or microbial co-inhabitants. Interestingly, MP1 and MP2 are

177    predicted to have retained the ability to produce spermidine, a known metabolic marker

178    of BV[40]. Despite the overall genomic reduction of MP1 and MP2, these vaginal phylotypes

179   have also gained functions specific to their clade. MP1 and MP2 specifically encode

180   virulence genes including variable tetracycline resistance genes (*i.e., tetM*, *tetO*, *tetW*)

181   and genes necessary for iron uptake (*i.e., tonB* and hemin uptake outer membrane

182   receptor). Iron sequestration is commonly a critical characteristic of pathogenic bacteria

183   and may be pertinent to the vaginal microbiome given the influx of available iron during

184   menses[41]. MP1 and MP2 genomes also encode multiple CRISPR-associated proteins,

185   which likely function to protect these bacteria from foreign genetic elements[42].

186   **Predicted functional divergence of MP1 and MP2**

187   MP1 and MP2 also possess unique predicted metabolic functions, indicative of

188   their divergence. While genomes of both phylotypes encode the majority of genes

189   required for glycolysis, MP2 genomes lack hexokinase. The absence of this gene

190   suggests that MP2 strains cannot use glucose as a carbon source.  MP1 genomes are

191   predicted to lack adenosine deaminase (ADA), an enzyme involved in the adenine

192   salvage pathway. In contrast, MP2 genomes retain ADA but lack the gene encoding

193   cytidine deaminase, which functions in the recycling of cytosine bases. These differential

194   salvage strategies are intriguing given that MP1 genomes have markedly higher GC

195   content than MP2 genomes. The phylotypes also differ in their ability to synthesize amino

196   acids. MP2 genomes are incapable of synthesizing leucine and tryptophan, while MP1

197   genomes lack the ability to interconvert serine and cysteine. Production of aromatic amino

198   acids including tryptophan is energetically expensive[43]. Thus, the loss of tryptophan

199   synthesis genes in MP2 is an example of energetically favorable genome reduction in this

200   host-associated organism.

201   **MP1 and MP2 phylotypes have distinct clinical associations**

202    Given the distinct metabolic capacities of the MP1 and MP2 phylotypes, we

203    examined the self-survey and clinical data associated with the Vaginal Human

204    Microbiome Project (VaHMP) to investigate their individual roles in reproductive health[27].

205    We first examined demographic and clinical associations with vaginal MP1 and MP2

206    carriage in a cohort of 3,091 non-pregnant women. In this cohort, 27% of women

207    (845/3091) carried MP1 only, 5% (163/3091) carried MP2 only, 6% (182/3091) carried

208    both phylotypes, and 62% (1901/3,091) carried neither phylotype. Compared to the

209    average alpha diversity (*i.e.,* inverse Simpson's index) of samples containing neither of

210    the two phylotypes (1.37), alpha diversity was increased in samples containing MP1 only

211    (1.79), MP2 only (3.47) and both phylotypes (3.37) (Fig. 4). Notably, vaginal microbiome

212    communities containing MP2 exhibited an almost two-fold increase in alpha diversity

213    compared with MP1 alone.

214    Associations with demographics were determined using a generalized linear

215    model. Both phylotypes were associated with African-ancestry (MP1: q= 3.00e-31, MP2:

216    q= 1.10e-21, with FDR correction) and a self-reported annual household income of less

217    than 20k (MP1: q= 2.23e-18, MP2: q= 3.31e-18 with FDR correction) (Supplementary

218    Table 7). Fethers *et. al* previously reported that MP1 was associated with women who

219    have sex with women (WSW)[44]. WSW experience higher rates of BV than women who

220    do not have sex with women[45]. Thus, we examined the association of both *Megasphaera*

221    phylotypes with WSW. Although 44% (38/86) of women who reported a current female

222    partner were MP1 positive and there was a positive association between WSW and MP1

223    carriage (q= 0.075 with FDR correction), it did not reach the threshold for significance

224    (p<0.05). Using the general linearized model (GLM), the race/ethnicity field was identified

225    as a significant covariate with WSW. In stratified analyses, we found that among women

226    who did not report African ancestry, there was a strong association between WSW (N=25)

227    and MP1 (q= 0.0012 with FDR correction), but that among women reporting African

228    ancestry, WSW (N=61) was not significantly associated with MP1 (q= 0.846 with FDR

229    correction). The majority of participants not reporting African ancestry self-identified as

230    Caucasian (68%). This finding highlights the need for precision medicine approaches that

231    account for the contribution of individual environmental and genetic factors and their

232    interactions to fully understand the contributions that shape vaginal microbiome

233    composition and impact risk for adverse reproductive health outcomes.

234         To assess the association of these two phylotypes with three common vaginal

235    infections (*i.e.,* bacterial vaginosis, candidiasis and trichomoniasis) we performed a

236    relative risk analysis. We observed that while both MP1 and MP2 were associated with

237    an increased risk for BV (MP1: 4.57, 95% CI [3.76,5.55], MP2: 2.19, 95% CI [1.79-2.69]),

238    MP1 is associated with a higher risk for this condition (Table 1). In contrast, MP2 was

239    associated with an increased risk for trichomoniasis (4.84, 95% CI [3.06-7.64]), whereas

240    MP1 had no association (0.96, 95% CI [0.59-1.56]). Using the GLM approach, MP1 and

241    MP2 strains were both associated with self-reported vaginal odor (MP1: q= 5.39e-18 ,

242    MP2: q= 1.36e-10 with FDR correction) and vaginal discharge (MP1: q= 1.40e-17, MP2:

243    q= 4.64e-7 with FDR correction).  Both phylotypes were also associated with clinician-

244    diagnosed elevated vaginal pH (>4.5) (MP1: q= 3.56e-34, MP2: q= 7.29e-12 with FDR

245    correction) consistent with previous reports. Carriage of MP1 and MP2 were also

246    associated with having more than 10 lifetime sexual partners (MP1: q= 0.00037, MP2: q=

11

247    4.65e-5 with FDR correction) and having more than one sexual partner in the past month

248    (MP1: q= 0.0002, MP2: q= 2.24e-5 with FDR correction).

249    **MP1 and MP2 in Pregnancy**

250    Recent studies have shown that the vaginal microbiome in pregnancy is

251    associated with decreased alpha diversity and dominance of protective *Lactobacillus*

252    species[46–49]. Similarly, BV-associated organisms have been shown to be less prevalent

253    in pregnant women [27,50]. Thus, not surprisingly in a case-matched cohort of 779 pregnant

254    and 779 non-pregnant women from the VaHMP study, we found that MP2 was

255    significantly decreased in pregnancy (q< 0.05, Mann-Whitney U test with FDR correction)

256    (Fig. 5). This finding is in agreement with previous work demonstrating that BV organisms

257    are often less prevalent in pregnancy[27,50]. In contrast, MP1 was not significantly excluded

258    in the pregnant cohort (q= 0.596, Mann-Whitney U test with FDR correction). MP1 has

259    been previously associated with risk for preterm birth[12,23,24]; additional studies will be

260    necessary to determine whether the ability of MP1 to persist throughout gestation has

261    implications for complications in pregnancy.

262    To determine whether the two vaginal phylotypes were functionally active in

263    pregnancy, we analyzed metatranscriptomic data from 57 samples collected from 52

264    pregnant women who delivered at term as a part of the case-control Preterm Birth cohort

265    from the Multi-'Omic Microbiome Study – Pregnancy Initiative (MOMS-PI)[23]. This is a

266    reanalysis of a subset of an existing dataset previously published in 2019[23,50]. In this

267    cohort, 43 samples contained transcripts assigned to MP1 while only one sample

268    contained transcripts assigned to MP2 (Supplementary Table 8, Supplementary Table 9),

269    consistent with our observation that MP2 seems to be less prevalent in pregnancy while

12

270    MP1 is maintained. Because MP2 was only detected in a single sample, we will focus on

271    the findings pertaining to MP1 here. The data showed that *in vivo* in the vaginal

272    environment, MP1 strains transcribed genes from 34 unique pathways. Notably, MP1

273    strains transcribed genes involved in butyrate production, which has previously been

274    associated with BV[40].

275         For this cohort (N=57), we also had paired 16S rDNA profiles and metagenome

276    sequencing profiles generated as a part of a previous study[23]. In these paired data, we

277    observed that the 16S rDNA relative abundance measures for MP1 were strongly

278    correlated to their paired metagenomic relative abundance measures ($\rho$=0.92,

279    Spearman's rank correlation). This finding supports the use of 16S rDNA profiles in lieu

280    of metagenomic sequencing data to estimate the relative abundance of MP1 in these

281    cohorts. The correlation of MP1 metagenomic relative abundance measures to their

282    paired metatranscriptomic relative abundance measures was also significant ($\rho$=0.91,

283    Spearman's rank correlation). Intriguingly, the relative abundance measures of the

284    transcripts assigned to MP1 were greater than the observed relative abundance

285    measures in the paired metagenomic dataset (p= 2.95e-05, Mann-Whitney U test) (Fig.

286    6). This suggests that MP1 is highly transcriptionally active in these samples and makes

287    up a greater proportion of the transcripts than would be predicted based upon the

288    metagenomic data alone. Taken together, the above analyses demonstrate that MP1 is

289    maintained in pregnancy, in contrast to other BV-associated organisms, and is

290    transcriptionally active in a majority of pregnant women in our cohort. These observations

291    in combination with previous associations of MP1 with PPROM and spontaneous preterm

292    labor, highlights MP1 as an important target for future study[12,23,24].

**DISCUSSION**

In conclusion, our phylogenetic analyses suggest that MP1 and MP2 are evolutionarily divergent from other *Megasphaera* species as well as each other. While comprehensive biological and physiological assays of MP1 and MP2 isolates would be necessary, there is strong phylogenetic evidence that supports placement of MP1 and MP2 into a separate genus. Compared to other *Megasphaera*, both organisms exhibit loss of gut-specific metabolic pathways, acquisition of iron uptake pathways, and loss of genes involved in the biosynthesis of differential amino acids. These organisms also exhibit reduced genomes and lowered GC composition, indicative of a transition to a more host-dependent state[47] that seems to be a common feature of adaption to the vaginal environment[35,36]. Taken together these observations are suggestive of adaptation to the host and/or vaginal environment.

Several lines of evidence support the hypothesis that MP1 and MP2 have adapted from an ancestral gastrointestinal tract strain to colonize the vaginal niche: i) the similarity of MP1 and MP2 to human gastrointestinal *Megasphaera* species, ii) the ubiquity of *Megasphaera* in the GI tracts of humans and mammals[31,32], iii) the streamlined genomes of MP1 and MP2, a common feature of strains identified in the human vagina, and iv) the physical proximity of the rectum and vagina. Based on our observations, we hypothesize that these two phylotypes share a common ancestor, likely a colonizer of the gastrointestinal tract. Their evolutionary divergence is characterized by progressive gene loss and genome reduction, common features among host-dependent organisms. These changes may be indicative of host dependence and/or adaptation to the vaginal environment specifically.

316    MP1 and MP2 are evolutionarily divergent and functionally distinct from one

317    another as well, and these findings have important implications for the contributions of

318    these unique phylotypes to vaginal infections and pregnancy complications[3,12–14,17,24].

319    Several lines of evidence show differential associations of these two phylotypes with

320    clinical diagnoses and demographic factors. As expected, our analyses confirmed that

321    MP1 is tightly correlated with BV as diagnosed by Amsel's criteria in a cohort of 3,091

322    non-pregnant women of reproductive age. This result is consistent with numerous

323    previous studies that have demonstrated the strong association of MP1 with BV and led

324    to its use as a biomarker for the diagnosis of the condition[8/18/2020 5:45:00 PM]. MP2 was also

325    associated with BV (RR= 2.19, 95% C.I. (1.79-2.69)) in the cohort, but to a lesser extent

326    than MP1 (RR= 4.57, 95% C.I. (3.76-5.55)). This finding is consistent with previous

327    reports of the specificity and sensitivity of MP1 and MP2 for BV diagnosis. While MP1

328    and MP2 have both been reported to have high specificity for BV ranging from 88.5%-

329    98.1% for MP1 and 98.9-100% for MP2 as diagnosed by Amsel's criteria, Nugent score

330    or a combination of both diagnostic measures, the sensitivity of MP2 (6.9%-31.0%) has

331    been reported to be significantly lower than that of MP1 (68.4%-95.1%)[16,51]. In the current

332    study, we observed an overall prevalence of 33.2% (n=1027) for MP1 and 11.2% (n=345)

333    for MP2 among non-pregnant women of reproductive age. However, our results also

334    suggest that the two major *Megasphaera* phylotypes may be associated with different

335    subtypes of vaginal dysbiosis.

336    In the cohort of 3,091 non-pregnant women of reproductive age, MP2 was strongly

337    associated with trichomoniasis whereas MP1 was not associated with the condition. To

338    our knowledge, the association of MP2 with trichomoniasis was first described by Martin

15

339   *et al.* in 2013, and our study confirms and extends this observation[17]. Martin *et al.* also

340   highlighted an observation from a 1992 study suggesting that *Trichomonas vaginalis*

341   infection was associated with intermediate flora as defined by Nugent score among

342   pregnant women[52]. Together, these findings highlight the need to distinguish between

343   related taxa, such as MP1 and MP2, in microbiome analyses in order to accurately define

344   the functionally relevant subtypes of vaginal dysbiosis and how they contribute to adverse

345   reproductive health outcomes.

346         In the current study, MP1 and MP2 were both more prevalent among women who

347   reported African ancestry (MP1 AA: 41.1%, MP1 Non-AA: 19.8%, MP2 AA: 15.9%, MP2

348   Non-AA: 3.1%), consistent with several previous reports including an analysis of the first

349   1,686 women enrolled in the VaHMP cohort[1,27,53]. The association of MP1 and MP2 with

350   African ancestry is consistent with the increased incidence of BV among women with

351   African ancestry[53–55]. In a recent study of 33 white women and 16 black women who were

352   BV negative as assayed by both Amsel's and Nugent's criteria, Beamer *et al.* did not

353   detect significant differences in the colonization and density of a number of bacterial

354   species assayed by cultivation and molecular methods[56]. Notably, organisms such as

355   MP1 are rarely observed in women with low Nugent scores; MP1 was identified in 2/33

356   (6.1%) white women and 2/16 black women (12.5%) by qPCR. While it is less refined

357   than new molecular methods for assaying the vaginal environment, Nugent score, which

358   calculates a score for BV based on the presence of bacterial morphotypes as assayed by

359   microscopy, is still a direct measure of microbial composition. By excluding individuals

360   with higher Nugent score, a significant proportion of women of African ancestry may be

361   excluded from the study. In this current analysis of the VaHMP cohort, race and ethnicity

16

362    were tightly correlated with a number of covariates including measures of socioeconomic

363    status such as education and annual income. Additional studies will be needed to further

364    define the contributions of both genetic and environmental factors that shape vaginal

365    microbiome composition[55,57,58].

366        In the current study, MP1 and MP2 were also found to be associated with

367    increased alpha diversity of the microbiome profile and elevated vaginal pH (>4.5), which

368    is one of Amsel's criteria for BV and is consistent with previous studies linking these

369    *Megasphaera* phylotypes to BV. Interestingly, the vaginal microbiome profile of women

370    who carried MP2 exhibited higher alpha diversity compared to the vaginal microbiome of

371    women who carried MP1 alone. MP2 exhibits greater genome reduction than MP1, likely

372    making it more reliant on other microbial species and/or host factors. This genomic

373    reduction may account for why MP2 is less prevalent in the overall population and specific

374    to a more diverse dysbiotic state.

375        MP1 is prevalent in the vaginal environment and has been associated with preterm

376    birth in several recent studies, marking it as a taxon of interest[23–25,49,50]. Our current study

377    suggests that MP1 levels are similar among pregnant and non-pregnant women, unlike

378    many other BV-associated vaginal taxa which seem to be excluded during the gestational

379    shaping of the vaginal microbiome[25,59,60]. MP1 is highly prevalent in the VaHMP cohort,

380    colonizing 33.2% of women in the study. These findings suggest that this highly prevalent

381    organism colonizes the vaginal environment and remains present and transcriptionally

382    active during pregnancy. Mitchell *et al.* observed MP1 in the upper genital tract (UGT) of

383    women undergoing hysterectomy suggesting that this organism is likely capable of

384    ascending into the UGT.  This capability combined with the ability of MP1 to maintain

17

385  colonization during pregnancy suggests that this organism is a candidate for future

386  studies investigating the proposed model where ascending infection of vaginal organisms

387  contributes to in preterm labor and/or birth. *Megasphaera* has also been associated with

388  low vitamin D levels[61] highlighting a possible link between vaginal microbiome signatures

389  and host state[62]. Identifying mechanisms that permit this organism to pervade the

390  changing vaginal environment associated with the progression of pregnancy may possibly

391  lead to the development of more effective preventative therapeutics targeting microbe-

392  related preterm labor and delivery. This study also highlights the need for continued

393  exploration of mechanisms of microbial evolution in the human microbiome.

394  Understanding the processes that underlie adaptation to specific human host-associated

395  environments will inform strategies for modulating the microbiome to prevent disease and

396  promote human microbial health.

397

398  **AUTHOR CONTRIBUTIONS**

399  A.L.G. conducted all experiments and analyses and drafted the manuscript. N.R.J.

400  contributed to clinical association analyses and manuscript preparation. S.B. contributed

401  to comparative genomic analyses. V.N.K. performed genome assembly and initial

402  genome annotation. J.P.B. and D.J.E. provided support for statistical analyses. J.F.S.

403  provided clinically relevant interpretation of results. K.K.J. oversaw cultivation of isolates

404  and provided interpretation of results. M.G.S. oversaw genome sequencing and provided

405  interpretation of phylogenetic data. G.A.B. contributed to interpretation of the results, and

406  G.A.B., K.K.J., J.F.S. and J.M.F. serve as the executive leadership team and planned

407  and directed the overall VaHMP and MOMS-PI studies. The VMC provided infrastructure

18

408    and data for the study. J.M.F. supervised this study and led the overall direction and

409    planning. A.L.G. and J.M.F designed the study and wrote the manuscript with

410    contributions from all other authors.

411

412    **ACKNOWLEDGMENTS**

430

431 **COMPETING INTERESTS**

432 Some authors (Brooks, Edwards, Strauss, Jefferson, Serrano, Buck & Fettweis) are co-

433 inventors on a pending patent for a preterm birth diagnostic signature. Strauss serves

434 as a Member on the Scientific Advisory Board of Prescient Medicine.

435

436 **MATERIALS AND METHODS**

437 **Cultivation of MP1 and MP2**

438 Using anaerobic technique, we cultivated, isolated and sequenced the genomes

439 of one isolate of *Megasphaera* phylotype 1 (MP1, strain M1-70) and two isolates of

440 *Megasphaera* phylotype 2 (MP2, strains M2-4 and M2-8) from frozen glycerol stocks of

441 vaginal swab samples collected through the Vaginal Human Microbiome Project

442 (VaHMP)[63]. One mid-vaginal swab from each participant was used to inoculate 1.0mL of

443 supplemented brain-heart infusion (sBHI) culture media with an added cryo-protectant

444 (20% glycerol) and stored at -80°C (Supplementary Table 11). Frozen vaginal culture

445 samples were targeted for cultivation based on the presence and high relative abundance

446 of bacterial targets of interest. These samples were identified using 16S rRNA gene

447 based vaginal microbiome profiles generated for each participant. A scraping of the frozen

448 vaginal culture media from the selected targets was used to inoculate agar plates for

449 bacterial culture. Scrapings were plated on both ThermoScientific Remel Chocolate agar

450 (lysed blood agar) and ThermoScientific Remel Brucella Blood agar (5% sheep's blood)

451 at four dilutions: 1:10, 1:100, 1:1000 and 1:10000. Plates were stored at 37°C for 24-48

452 hours. The plates were enclosed in three nested Ziploc bags along with a Mitsubishi

453 Anaeropack-Anaero packet to simulate anaerobic conditions. Individual colonies were

20

454 selected for growth and purification from the dilution plates based on colony morphology

455 and differential growth characteristics. After re-streaking for visibly pure colonies, the

456 isolates were taxonomically identified by colony PCR amplification of the full 16S rRNA

457 gene using universal 16S primers[64,65]. Amplicons were purified using the Qiagen

458 QIAquick PCR Purification Kit and sequenced using the Applied Biosystems 3730 DNA

459 Analyzer. Colonies that were identified as bacterial targets of interest and exhibited no

460 evidence of contamination were selected for extraction of genomic DNA. A single colony

461 inoculum was added to 5mL of sBHI in a 15mL falcon tube. Tubes were loosely capped

462 to allow gas exchange and stored in a rack at 37°C for 24-48 hours in three nested Ziploc

463 bags containing a Mitsubishi Anaeropack-Anaero. The DNA was then extracted using the

464 Qiagen DNeasy Blood & Tissue Kit and quantified using the Nanodrop 2000

465 spectrophotometer. Frozen stocks for MP1 and MP2 isolates were not recoverable.

466 **Genome Sequencing and Assembly**

467 Purified genomic DNA from the single MP1 isolate was sequenced using the

468 Roche 454 GS FLX Titanium platform. The resulting reads were trimmed for quality and

469 assembled using Newbler v2.8[66]. Purified genomic DNA derived from the two MP2

470 isolates M2-4 and M2-8 were sequenced using the Illumina MiSeq platform and the

471 resulting reads were trimmed for quality and assembled using Newbler v2.8, CLCBio and

472 SPAdes[66–68]. These three assemblies were merged using CISA to produce the most

473 complete and accurate contigs[69].

474 **Structural Genomic Analysis**

475 Genomic synteny was analyzed between genome representatives of MP1 and

476 MP2 and other host-associated *Megasphaera* species. This analysis was performed at

477    the both the protein and nucleic acid level. Nucleic acid-based synteny analyses were

478    performed using NUCmer while amino acid-based synteny analyses were performed

479    using PROmer. Both NUCmer and PROmer are available as a part of the MUMmer 3.0

480    package[70]. Synteny plots were created using gnuplot from the gnuplot 4.2 package and

481    MUMmerplot, which is also available as a part of the MUMmer 3.0 package[71]. Genomic

482    GC composition was determined using in-house scripts. Codon usage within the

483    genomes was calculated using cusp, a program included in the EMBOSS Tools package

484    available through EMBL-EBI[72]. Comparative analyses of basic genome statistics

485    including genome size, predicted number of proteins and GC composition were

486    performed using a Kruskal-Wallis test. This was performed using the kruskal.test function

487    in R. All calculated p values were adjusted using the FDR correction in R using the

488    p.adjust function. Resulting corrected q values are reported in the Results.

489    **Measures of Genomic Similarity**

490         Analyses were performed utilizing three MP1 genomes, three MP2 genomes and

491    all publicly available *Megasphaera* and *Anaeroglobus* genomes at NCBI as of January 1,

492    2015 (Supplementary Table 10). One metagenomic *Megasphaera elsdenii* assembly was

493    excluded from the analysis due to variation in size and gene content from other deposited

494    *M. elsdenii* genomes. One representative of MP1 (Veillonellaceae bacterium DNF00751)

495    and two representatives of MP2 (*Megasphaera* genomosp. 2, Veillonellaceae bacterium

496    KA00182) were deposited after analyses were complete. ANI values suggest that they

497    are similar in genomic content to the genome representatives analyzed in this study.

498    Veillonellaceae bacterium DNF00751 had ANI values ranging from 96.5-98.6% compared

499    to the three MP1 genomes utilized in our analysis. *Megasphaera* genomosp. 2 had ANI

22

500   values ranging from 98.5-99.0% and Veillonellaceae bacterium KA00182 had ANI values

501   ranging from 98.6-99.0% to the three MP2 genomes utilized in our analyses.

502   To assess genomic similarity using the entire nucleotide content of the genomes,

503   a pairwise calculation of the average nucleotide identity was performed using a publicly

504   available script (https://github.com/chjp/ANI). 16S ribosomal RNA gene sequences are

505   commonly used to distinguish bacterial species and establish evolutionary relatedness[73].

506   16S rRNA gene sequences were identified and extracted from genomes using

507   RNAmmer[74]. Sequence similarity of the 16S rRNA genes was determined using the blastn

508   algorithm[75]. In order to delineate genus boundaries, pairwise Percentage of Conserved

509   Proteins (POCP) values were calculated using in-house scripts developed based on the

510   methods described in Qin et al., 2014[28].

511   **CSI and CSP detection**

512   Conserved Signature Proteins (CSPs) and genomic regions containing Conserved

513   Signature Indels (CSIs) were identified using BLAST [75]. Genomic regions containing CSIs

514   were aligned using MUSCLE and visualized using Jalview[76,77]. This analysis was based

515   on work performed by Campbell *et al.* identifying CSPs and CSIs indicative of the

516   placement of certain taxa within the class Negativicutes[30].

517   **Genome Annotation and Metabolic Reconstruction**

518   Genomes were annotated using both an in-house annotation pipeline and RAST[77],

519   a web-based genome visualization, annotation and metabolic reconstruction tool

520   provided by NMPDR[78]. As a part of the in-house Genome Annotation Pipeline, the

521   following programs were used. Genes were called using both Glimmer3 and

522   GeneMarkS[79,80]. Ribosomal RNA genes were identified and extracted from genomes

23

523    using RNAmmer[74]. Genes encoding tRNAs were identified in genomes using tRNAScan-

524    SE[81]. Orthologous genes were detected using rpsblast in conjunction with Pfam and COG

525    databases[75,82–85]. Predicted gene functions were annotated using blastx and the nr

526    database at NCBI[75]. Metabolic reconstruction was performed using ASGARD[85] and visual

527    representations of predicted variation within metabolic pathways were generated using

528    the program color-maps[86]. To determine genes lost in MP1 and MP2, genes specific to

529    MP1 and MP2 and genes that can be used to distinguish the two phylotypes, RAST

530    annotation was utilized. Findings were verified by comparing RAST results to the Genome

531    Annotation Pipeline Glimmer3 and GeneMarkS gene calls. Further verification was

532    performed using the tblastn algorithm to compare known annotated protein sequences

533    available through NCBI to the raw genomic contigs[75,82,85].

**Phylogenetic Analysis**

535        To perform a phylogenetic reconstruction of the 16S rRNA gene, 16S rRNA

536    sequences were identified and extracted from the genomes using RNAmmer[74]. The

537    extracted 16S rRNA gene sequences were aligned using MUSCLE[76]. The resulting

538    alignment file was converted to phylip format using a web-based tool for DNA and protein

539    file format conversion, **AL**ignment **T**ransformation **E**nvi**R**onment or ALTER[87]. RAxML-

540    HPC was used to perform a rapid bootstrap analysis using 1,000 bootstraps and search

541    for the best scoring maximum likelihood tree using the gamma model of heterogeneity[88].

542    To create a phylogenetic reconstruction of all Negativicutes class genomes, 145

543    orthologous genes were used. OrthoDB, an online database for orthologous groups was

544    used to determine which orthologous genes were conserved at the family level

545    (Veillonellaceae)[89]. These genes were verified using reciprocal blast and extracted from

24

546    the six MP1 and MP2 genomes as well as from all publicly available genomes classified

547    to the class Negativicutes at NCBI as of January 1, 2015. *Clostridium botulinum* A strain

548    Hall was selected as the outgroup. This species was chosen due to its classification in

549    the same phylum (Firmicutes) but different class (Clostridia versus Negativicutes) as

550    compared to the Negativicutes genomes. Each orthologous gene was separately aligned

551    using MUSCLE, a program within the EMBOSS Tools package available through EMBL-

552    EBI[72,76]. Alignments were visually examined and those with large gaps or likely errors

553    were discarded. Sequences from all orthologs were concatenated together to form one

554    large informative sequence. Concatenated sequences were then pruned for informative

555    regions using Gblocks[90]. The resulting sequences were converted from pir to phylip

556    format using the web tool, ALTER[87]. RAxML-HPC was used to perform a rapid bootstrap

557    analysis using 100 bootstraps and search for the best scoring maximum likelihood tree

558    using optimization of substitution rates, the gamma model of heterogeneity and the WAG

559    amino acid substitution matrix [88]. Aesthetic changes to the tree were made using

560    TreeDyn[91].

561    **Participant Recruitment and Informed Consent**

562        We used samples and data from two existing cohorts for this study, The Vaginal

563    Human Microbiome Project (VaHMP) and the Multi-Omic Microbiome Study–Pregnancy

564    Initiative (MOMS-PI), reviewed and approved by the Institutional Review Board at Virginia

565    Commonwealth University (IRB #HM12169, IRB #HM15527). Samples and data are

566    maintained in the Research Alliance for Microbiome Science (RAMS) Registry at Virginia

567    Commonwealth University (IRB #HM15528). The study was performed with compliance

568    to all relevant ethical regulations. Written informed consent was obtained for all

569  participants and parental permission and assent was obtained for participating minors at

570  least 15 years of age.

**Sample Collection, Vaginal Microbiome Profiling and Analysis**

572  Samples collected as part of the Vaginal Human Microbiome Project (VaHMP) at

573  Virginia Commonwealth University were used for this study as previously described[63].

574  Briefly, mid-vaginal wall swab samples were collected and DNA was extracted from the

575  swabs using the MoBio Powersoil DNA Isolation Kit. DNA samples were randomized to

576  avoid batch effects and the V1-V3 region of the 16S rRNA gene was amplified using

577  polymerase chain reaction (PCR) and universal primers (Supplementary Table 12)[64,65].

578  The amplified 16S rDNA fragments were sequenced using the Roche 454 GS FLX

579  Titanium platform. Sequences were classified using both the Ribosomal Database Project

580  (RDP) classifier and the in-house STIRRUPS (Species-level Taxon Identification of rDNA

581  Reads using a USEARCH Pipeline Strategy) classifier to achieve species-level

582  classification (version 10-18-17)[92,93]. Samples that yielded less that 5,000 reads were

583  excluded from analysis.

584  Taxa were determined to be present if they comprised at least 0.1% of the vaginal

585  microbiome profile of a given sample. Demographics and health history data was self-

586  reported by the participants. Associations were calculated based on the presence or

587  absence of a taxon of interest (threshold of 0.1% of total reads) in combination with given

588  demographic or clinical data. Statistical significance was calculated using a generalized

589  linear model using logistic regression as implemented in the 'glm' function in R. All

590  calculated p values were corrected for multiple testing using the FDR correction method.

26

591    This was performed in R using the p.adjust function. Adjusted q values are reported in

592    the Results.

593    **Alpha Diversity Measures**

594          16S rDNA-based vaginal microbiome profiles from the Vaginal Human Microbiome

595    Project (VaHMP) outpatient cohort of non-pregnant subjects was used for this analysis

596    (n=3091). Alpha diversity for each microbiome profile was calculated using relative

597    proportion data, renormalized to exclude unclassified reads (below 97% threshold).

598    Inverse Simpson's Index was used as the measure of alpha diversity. This metric was

599    calculated using the R package 'vegan.' Average Inverse Simpson's Index alpha diversity

600    measures were generated for four subsets of vaginal microbiome data i) samples

601    containing neither phylotype ii) samples containing only MP1 iii) samples containing only

602    MP2 and iv) samples containing both MP1 and MP2. Presence of a phylotype was

603    denoted by a relative abundance of greater than or equal to 0.1% of the vaginal

604    microbiome profile. Statistical significance was calculated using a two-tailed Student's T-

605    test with a significance level of 0.05.

606    **Relative Risk**

607          The non-pregnant, outpatient VaHMP cohort was used for this analysis (n=3091).

608    Samples met the threshold of at least 5,000 reads. Vaginal infection status was

609    determined based on clinician diagnosis at time of visit. Relative risk values and their

610    corresponding 95% confidence interval values were calculated based on the standard

611    relative risk formula. Relative Risk = (A/A+B) / (C/C+D) where A represents the number

612    of samples where the taxon is present and the participant is diagnosed with the disease,

613    B represents the number of samples where the taxon is present but the participant is not

27

614     diagnosed with the disease, C represents the number of samples where the taxon is

615     absent but the participant is diagnosed with the disease and D represents the number of

616     samples where the taxon is not present and the participant is not diagnosed with the

617     disease.

618     **Pregnancy Analysis**

619     A case-matched cohort was used for this analysis. A cohort of 779 pregnant

620     women was case-matched 1:1 based on ethnicity, age and income to 779 non-pregnant

621     controls. Using the R package 'wilcox', we performed a Mann-Whitney U test on all

622     vaginal microbial taxa present in at least 5% of samples that comprise at least 0.1%

623     relative proportion of the microbiome profile. We utilized the R function 'p.adjust' to correct

624     for multiple testing using the FDR correction. Results for three *Lactobacillus* species,

625     MP1, MP2 and select associated organisms associated with dysbiosis are shown[94].

626     **Transcriptomic Analyses**

627     The Multi-Omic Microbiome Study-Pregnancy Initiative (MOMS-PI) Preterm Birth

628     cohort was utilized for this analysis[95]. This cohort consists of several hundred thousand

629     samples collected from pregnant women throughout and after their pregnancies. For

630     meta-transcriptomics, we collected a mid-vaginal swab from each participant and pre-

631     processed the sample within an hour of collection by inserting the swab into RNAlater®

632     (Qiagen). These swabs were then processed using MoBio Power Microbiome RNA

633     Isolation kit as described by the manufacturer. Total RNA was depleted of human and

634     microbial rRNA using the Epicentre/Illumina Ribo-Zero Magnetic Epidemiology Kit as

635     described by the manufacturer. Enriched messenger RNA was prepared for sequencing

636     by constructing cDNA libraries using the KAPA Biosystems KAPA RNA HyperPrep Kit.

637　Indexed cDNA libraries were pooled in equimolar amounts and sequenced on the Illumina

638　HiSeq 4000 instrument running 4 multiplexed samples per lane with an average yield of

639　~100 Gb/lane, sufficient to provide >100X coverage of the expression profiles of the most

640　abundant 15-20 taxa in a sample. Raw sequence data was demultiplexed into sample-

641　specific fastq files using *bcl2fastq* conversion software from Illumina. Adapter residues

642　were trimmed from both 5' and 3' end of the reads using Adapter Removal tool v2.1.3.

643　The sequences were trimmed for quality using *meeptools*[96], retaining reads with minimum

644　read length of 70b and *meep* (maximum expected error) quality score less than 1. Human

645　reads were identified and removed from each sample by aligning the reads to hg19 build

646　of the human genome using the BWA aligner[97]. Transcripts were classified using

647　HUMAnN2[98,99] and shortBRED[100]. Transcripts assigned to either MP1 or MP2 were

648　analyzed for this study.

649　**Data Availability**

650　　The genomes of *Megasphaera* phylotype 1 (MP1, strain M1-70), *Megasphaera*

651　phylotype 2 (MP2, strain M2-4) and *Megasphaera* phylotype 2 (MP2, strain M2-8) have

652　been submitted to DDBJ/ENA/GenBank under accession numbers PTJT00000000,

653　PTJU00000000 and PTJV00000000 respectively. The versions described in this paper

654　are versions PTJT01000000, PTJU01000000 and PTJV01000000. Data from the VaHMP

655　has been deposited under dbGAP Study Accession phs000256.v3.p2. Raw

656　metatranscriptomic sequences from the MOMS-PI project are available at NCBI's

657　controlled-access dbGaP (Study Accession: phs001523.v1.p1). Access to additional

658　fields can be requested through the RAMS Registry (https://ramsregistry.vcu.edu).

659　**Code availability**

29

660 Custom code for GC composition and Percentage of Conserved Protein (POCP)

661 calculations is available at https://github.com/Vaginal-Microbiome-Consortium.

662
663 **BIBLIOGRAPHY**
664

665 1. Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc. Natl. Acad.*

666 *Sci. U. S. A.* **108 Suppl 1**, 4680–4687 (2011).

667 2. Romero, R. *et al.* The role of infection in preterm labour and delivery. *Paediatr.*

668 *Perinat. Epidemiol.* **15 Suppl 2**, 41–56 (2001).

669 3. Fethers, K. *et al.* Bacterial vaginosis (BV) candidate bacteria: associations with BV

670 and behavioural practices in sexually-experienced and inexperienced women. *PloS*

671 *One* **7**, e30633 (2012).

672 4. van de Wijgert, J. H. H. M. The vaginal microbiome and sexually transmitted

673 infections are interlinked: Consequences for treatment and prevention. *PLoS Med.*

674 **14**, e1002478 (2017).

675 5. Dominguez-Bello, M. G. *et al.* Delivery mode shapes the acquisition and structure of

676 the initial microbiota across multiple body habitats in newborns. *Proc. Natl. Acad.*

677 *Sci. U. S. A.* **107**, 11971–11975 (2010).

678 6. van Best, N., Hornef, M. W., Savelkoul, P. H. M. & Penders, J. On the origin of

679 species: Factors shaping the establishment of infant's gut microbiota. *Birth Defects*

680 *Res. Part C Embryo Today Rev.* **105**, 240–251 (2015).

681 7. Tachedjian, G., Aldunate, M., Bradshaw, C. S. & Cone, R. A. The role of lactic acid

682 production by probiotic Lactobacillus species in vaginal health. *Res. Microbiol.* **168**,

683 782–792 (2017).

684    8.  Peebles, K., Velloza, J., Balkus, J. E., McClelland, R. S. & Barnabas, R. V. High

685        Global Burden and Costs of Bacterial Vaginosis: A Systematic Review and Meta-

686        Analysis. *Sex. Transm. Dis.* **46**, 304–311 (2019).

687    9.  Cohen, C. R. *et al.* Bacterial Vaginosis Associated with Increased Risk of Female-to-

688        Male HIV-1 Transmission: A Prospective Cohort Analysis among African Couples.

689        *PLoS Med* **9**, e1001251 (2012).

690    10. Hillier, S. L. *et al.* The role of bacterial vaginosis and vaginal bacteria in amniotic

691        fluid infection in women in preterm labor with intact fetal membranes. *Clin. Infect.*

692        *Dis. Off. Publ. Infect. Dis. Soc. Am.* **20 Suppl 2**, S276-278 (1995).

693    11. Kenyon, C., Colebunders, R. & Crucitti, T. The global epidemiology of bacterial

694        vaginosis: a systematic review. *Am. J. Obstet. Gynecol.* **209**, 505–523 (2013).

695    12. Nelson, D. B. *et al.* Early pregnancy changes in bacterial vaginosis-associated

696        bacteria and preterm delivery. *Paediatr. Perinat. Epidemiol.* **28**, 88–96 (2014).

697    13. Zozaya-Hinchliffe, M., Martin, D. H. & Ferris, M. J. Prevalence and abundance of

698        uncultivated Megasphaera-like bacteria in the human vaginal environment. *Appl.*

699        *Environ. Microbiol.* **74**, 1656–1659 (2008).

700    14. Datcu, R. *et al.* Bacterial vaginosis diagnosed by analysis of first-void-urine

701        specimens. *J. Clin. Microbiol.* **52**, 218–225 (2014).

702    15. Lennard, K. *et al.* Microbial Composition Predicts Genital Tract Inflammation and

703        Persistent Bacterial Vaginosis in South African Adolescent Females. *Infect. Immun.*

704        **86**, (2018).

705  16. Fredricks, D. N., Fiedler, T. L., Thomas, K. K., Oakley, B. B. & Marrazzo, J. M.

706      Targeted PCR for detection of vaginal bacteria associated with bacterial vaginosis.

707      *J. Clin. Microbiol.* **45**, 3270–3276 (2007).

708  17. Martin, D. H. *et al.* Unique vaginal microbiota that includes an unknown

709      Mycoplasma-like organism is associated with Trichomonas vaginalis infection. *J.*

710      *Infect. Dis.* **207**, 1922–1931 (2013).

711  18. McClelland, R. S. *et al.* Evaluation of the association between the concentrations of

712      key vaginal bacteria and the increased risk of HIV acquisition in African women from

713      five cohorts: a nested case-control study. *Lancet Infect. Dis.* **18**, 554–564 (2018).

714  19. Sabo, M. C. *et al.* Associations between vaginal bacteria implicated in HIV

715      acquisition risk and proinflammatory cytokines and chemokines. *Sex. Transm.*

716      *Infect.* (2019) doi:10.1136/sextrans-2018-053949.

717  20. Zozaya, M. *et al.* Bacterial communities in penile skin, male urethra, and vaginas of

718      heterosexual couples with and without bacterial vaginosis. *Microbiome* **4**, (2016).

719  21. Nelson, D. E. *et al.* Bacterial communities of the coronal sulcus and distal urethra of

720      adolescent males. *PloS One* **7**, e36298 (2012).

721  22. Lamont, R. F. Advances in the Prevention of Infection-Related Preterm Birth. *Front.*

722      *Immunol.* **6**, (2015).

723  23. Fettweis, J. M. *et al.* The vaginal microbiome and preterm birth. *Nat. Med.* **25**, 1012–

724      1021 (2019).

725  24. Paramel Jayaprakash, T. *et al.* High Diversity and Variability in the Vaginal

726      Microbiome in Women following Preterm Premature Rupture of Membranes

727      (PPROM): A Prospective Cohort Study. *PloS One* **11**, e0166794 (2016).

728  25. Hočevar, K. *et al.* Vaginal Microbiome Signature Is Associated With Spontaneous

729      Preterm Delivery. *Front. Med.* **6**, 201 (2019).

730  26. Mitchell, C. M. *et al.* Colonization of the upper genital tract by vaginal bacterial

731      species in nonpregnant women. *Am. J. Obstet. Gynecol.* **212**, 611.e1–9 (2015).

732  27. Fettweis, J. M. *et al.* Differences in vaginal microbiome in African American women

733      versus women of European ancestry. *Microbiol. Read. Engl.* **160**, 2272–2282

734      (2014).

735  28. Qin, Q.-L. *et al.* A proposed genus boundary for the prokaryotes based on genomic

736      insights. *J. Bacteriol.* **196**, 2210–2215 (2014).

737  29. Carlier, J.-P. *et al.* Anaeroglobus geminatus gen. nov., sp. nov., a novel member of

738      the family Veillonellaceae. *Int. J. Syst. Evol. Microbiol.* **52**, 983–986 (2002).

739  30. Campbell, C., Adeolu, M. & Gupta, R. S. Genome-based taxonomic framework for

740      the class Negativicutes: division of the class Negativicutes into the orders

741      Selenomonadales emend., Acidaminococcales ord. nov. and Veillonellales ord. nov.

742      *Int. J. Syst. Evol. Microbiol.* **65**, 3203–3215 (2015).

743  31. Rogosa, M. Transfer of Peptostreptococcus elsdenii Gutierrez et al. to a New

744      Genus, Megasphaera [M. elsdenii (Gutierrez et al.) comb. nov.]. *Int. J. Syst.*

745      *Bacteriol.* **21**, 187–189 (1971).

746  32. Padmanabhan, R. *et al.* Non-contiguous finished genome sequence and description

747      of Megasphaera massiliensis sp. nov. *Stand. Genomic Sci.* **8**, 525–538 (2013).

748  33. Marchandin, H. *et al.* Phylogenetic analysis of some Sporomusa sub-branch

749      members isolated from human clinical specimens: description of Megasphaera

750      micronuciformis sp. nov. *Int. J. Syst. Evol. Microbiol.* **53**, 547–553 (2003).

751    34. Shetty, S. A., Marathe, N. P., Lanjekar, V., Ranade, D. & Shouche, Y. S.

752         Comparative genome analysis of Megasphaera sp. reveals niche specialization and

753         its potential role in the human gut. *PloS One* **8**, e79353 (2013).

754    35. Mendes-Soares, H., Suzuki, H., Hickey, R. J. & Forney, L. J. Comparative functional

755         genomics of Lactobacillus spp. reveals possible mechanisms for specialization of

756         vaginal lactobacilli to their environment. *J. Bacteriol.* **196**, 1458–1470 (2014).

757    36. Yeoman, C. J. *et al.* Comparative Genomics of Gardnerella vaginalis Strains

758         Reveals Substantial Differences in Metabolic and Virulence Potential. *PLoS One* **5**,.

759    37. Mende, D. R., Sunagawa, S., Zeller, G. & Bork, P. Accurate and universal

760         delineation of prokaryotic species. *Nat. Methods* **10**, 881–884 (2013).

761    38. Srinivasan, S. *et al.* Characterization of novel megasphaera species from the female

762         reproductive tract. *Am. J. Obstet. Gynecol.* **219**, 648–649 (2018).

763    39. Kim, M., Oh, H.-S., Park, S.-C. & Chun, J. Towards a taxonomic coherence between

764         average nucleotide identity and 16S rRNA gene sequence similarity for species

765         demarcation of prokaryotes. *Int. J. Syst. Evol. Microbiol.* **64**, 346–351 (2014).

766    40. Srinivasan, S. *et al.* Metabolic signatures of bacterial vaginosis. *mBio* **6**, (2015).

767    41. Cornelissen, C. N. & Hollander, A. TonB-Dependent Transporters Expressed by

768         Neisseria gonorrhoeae. *Front. Microbiol.* **2**, 117 (2011).

769    42. Koonin, E. V. & Makarova, K. S. Mobile Genetic Elements and Evolution of

770         CRISPR-Cas Systems: All the Way There and Back. *Genome Biol. Evol.* **9**, 2812–

771         2825 (2017).

772    43. Priya, V. K., Sarkar, S. & Sinha, S. Evolution of tryptophan biosynthetic pathway in

773         microbial genomes: a comparative genetic study. *Syst. Synth. Biol.* **8**, 59–72 (2014).

774    44. Fethers, K. *et al.* Bacterial vaginosis (BV) candidate bacteria: associations with BV

775        and behavioural practices in sexually-experienced and inexperienced women. *PloS*

776        *One* **7**, e30633 (2012).

777    45. Plummer, E. L. *et al.* Sexual practices have a significant impact on the vaginal

778        microbiota of women who have sex with women. *Sci. Rep.* **9**, 19749 (2019).

779    46. MacIntyre, D. A. *et al.* The vaginal microbiome during pregnancy and the

780        postpartum period in a European population. *Sci. Rep.* **5**, 8988 (2015).

781    47. Romero, R. *et al.* The composition and stability of the vaginal microbiota of normal

782        pregnant women is different from that of non-pregnant women. *Microbiome* **2**, 4

783        (2014).

784    48. Aagaard, K. *et al.* A metagenomic approach to characterization of the vaginal

785        microbiome signature in pregnancy. *PloS One* **7**, e36466 (2012).

786    49. Walther-António, M. R. S. *et al.* Pregnancy's stronghold on the vaginal microbiome.

787        *PloS One* **9**, e98514 (2014).

788    50. Serrano, M. G. *et al.* Racioethnic diversity in the dynamics of the vaginal microbiome

789        during pregnancy. *Nat. Med.* **25**, 1001–1011 (2019).

790    51. Hilbert, D. W. *et al.* Development and Validation of a Highly Accurate Quantitative

791        Real-Time PCR Assay for Diagnosis of Bacterial Vaginosis. *J. Clin. Microbiol.* **54**,

792        1017–1024 (2016).

793    52. Hillier, S. L., Krohn, M. A., Nugent, R. P. & Gibbs, R. S. Characteristics of three

794        vaginal flora patterns assessed by gram stain among pregnant women. Vaginal

795        Infections and Prematurity Study Group. *Am. J. Obstet. Gynecol.* **166**, 938–944

796        (1992).

797    53. Marrazzo, J. M. Interpreting the epidemiology and natural history of bacterial

798        vaginosis: are we still confused? *Anaerobe* **17**, 186–190 (2011).

799    54. Ness, R. B. *et al.* Can known risk factors explain racial differences in the occurrence

800        of bacterial vaginosis? *J. Natl. Med. Assoc.* **95**, 201–212 (2003).

801    55. Borgdorff, H. *et al.* The association between ethnicity and vaginal microbiota

802        composition in Amsterdam, the Netherlands. *PloS One* **12**, e0181135 (2017).

803    56. Beamer, M. A. *et al.* Bacterial species colonizing the vagina of healthy women are

804        not associated with race. *Anaerobe* **45**, 40–43 (2017).

805    57. Taylor, B. D. *et al.* Toll-like receptor variants and cervical Atopobium vaginae

806        infection in women with pelvic inflammatory disease. *Am. J. Reprod. Immunol. N. Y.*

807        *N 1989* **79**, (2018).

808    58. Murphy, K. & Mitchell, C. M. The Interplay of Host Immunity, Environment and the

809        Risk of Bacterial Vaginosis and Associated Reproductive Health Outcomes. *J.*

810        *Infect. Dis.* **214 Suppl 1**, S29-35 (2016).

811    59. Chu, D. M., Seferovic, M., Pace, R. M. & Aagaard, K. M. The microbiome in preterm

812        birth. *Best Pract. Res. Clin. Obstet. Gynaecol.* **52**, 103–113 (2018).

813    60. Brown, R. G. *et al.* Establishment of vaginal microbiota composition in early

814        pregnancy and its association with subsequent preterm prelabor rupture of the fetal

815        membranes. *Transl. Res. J. Lab. Clin. Med.* **207**, 30–43 (2019).

816    61. Jefferson, K. K. *et al.* Relationship between vitamin D status and the vaginal

817        microbiome during pregnancy. *J. Perinatol. Off. J. Calif. Perinat. Assoc.* **39**, 824–836

818        (2019).

819    62. Integrative HMP (iHMP) Research Network Consortium. The Integrative Human

820        Microbiome Project. *Nature* **569**, 641–648 (2019).

821    63. Fettweis, J. M., Serrano, M. G., Girerd, P. H., Jefferson, K. K. & Buck, G. A. A new

822        era of the vaginal microbiome: advances using next-generation sequencing. *Chem.*

823        *Biodivers.* **9**, 965–976 (2012).

824    64. Frank, J. A. *et al.* Critical evaluation of two primers commonly used for amplification

825        of bacterial 16S rRNA genes. *Appl. Environ. Microbiol.* **74**, 2461–2470 (2008).

826    65. Romero, R. *et al.* The vaginal microbiota of pregnant women who subsequently

827        have spontaneous preterm labor and delivery and those with a normal delivery at

828        term. *Microbiome* **2**, 18 (2014).

829    66. Miller, J. R., Koren, S. & Sutton, G. Assembly Algorithms for Next-Generation

830        Sequencing Data. *Genomics* **95**, 315–327 (2010).

831    67. Krämer, A., Green, J., Pollard, J. & Tugendreich, S. Causal analysis approaches in

832        Ingenuity Pathway Analysis. *Bioinforma. Oxf. Engl.* **30**, 523–530 (2014).

833    68. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its

834        Applications to Single-Cell Sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).

835    69. Lin, S.-H. & Liao, Y.-C. CISA: Contig Integrator for Sequence Assembly of Bacterial

836        Genomes. *PLoS ONE* **8**, (2013).

837    70. Delcher, A. L., Phillippy, A., Carlton, J. & Salzberg, S. L. Fast algorithms for large-

838        scale genome alignment and comparison. *Nucleic Acids Res.* **30**, 2478–2483

839        (2002).

840    71. gnuplot homepage. http://gnuplot.info/.

841     72. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open

842         Software Suite. *Trends Genet. TIG* **16**, 276–277 (2000).

843     73. Janda, J. M. & Abbott, S. L. 16S rRNA Gene Sequencing for Bacterial Identification

844         in the Diagnostic Laboratory: Pluses, Perils, and Pitfalls. *J. Clin. Microbiol.* **45**,

845         2761–2764 (2007).

846     74. Lagesen, K. *et al.* RNAmmer: consistent and rapid annotation of ribosomal RNA

847         genes. *Nucleic Acids Res.* **35**, 3100–3108 (2007).

848     75. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local

849         alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

850     76. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high

851         throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).

852     77. Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J.

853         Jalview Version 2--a multiple sequence alignment editor and analysis workbench.

854         *Bioinforma. Oxf. Engl.* **25**, 1189–1191 (2009).

855     78. Overbeek, R. *et al.* The SEED and the Rapid Annotation of microbial genomes using

856         Subsystems Technology (RAST). *Nucleic Acids Res.* **42**, D206-214 (2014).

857     79. Delcher, A. L., Harmon, D., Kasif, S., White, O. & Salzberg, S. L. Improved microbial

858         gene identification with GLIMMER. *Nucleic Acids Res.* **27**, 4636–4641 (1999).

859     80. Besemer, J., Lomsadze, A. & Borodovsky, M. GeneMarkS: a self-training method for

860         prediction of gene starts in microbial genomes. Implications for finding sequence

861         motifs in regulatory regions. *Nucleic Acids Res.* **29**, 2607–2618 (2001).

862     81. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of

863         transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).

864  82. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.

865      *Bioinforma. Oxf. Engl.* **26**, 2460–2461 (2010).

866  83. Finn, R. D. *et al.* The Pfam protein families database. *Nucleic Acids Res.* **38**, D211-

867      222 (2010).

868  84. Tatusov, R. L. *et al.* The COG database: an updated version includes eukaryotes.

869      *BMC Bioinformatics* **4**, 41 (2003).

870  85. Marchler-Bauer, A. *et al.* CDD: a database of conserved domain alignments with

871      links to domain three-dimensional structure. *Nucleic Acids Res.* **30**, 281–283 (2002).

872  86. Alves, J. M. P. & Buck, G. A. Automated System for Gene Annotation and Metabolic

873      Pathway Reconstruction Using General Sequence Databases. *Chem. Biodivers.* **4**,

874      2593–2602 (2007).

875  87. Glez-Peña, D., Gómez-Blanco, D., Reboiro-Jato, M., Fdez-Riverola, F. & Posada, D.

876      ALTER: program-oriented conversion of DNA and protein alignments. *Nucleic Acids*

877      *Res.* **38**, W14-18 (2010).

878  88. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis

879      of large phylogenies. *Bioinforma. Oxf. Engl.* **30**, 1312–1313 (2014).

880  89. Kriventseva, E. V. *et al.* OrthoDB v8: update of the hierarchical catalog of orthologs

881      and the underlying free software. *Nucleic Acids Res.* **43**, D250-256 (2015).

882  90. Castresana, J. Selection of conserved blocks from multiple alignments for their use

883      in phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540–552 (2000).

884  91. Chevenet, F., Brun, C., Bañuls, A.-L., Jacq, B. & Christen, R. TreeDyn: towards

885      dynamic graphics and annotations for analyses of trees. *BMC Bioinformatics* **7**, 439

886      (2006).

887   92. Fettweis, J. M. *et al.* Species-level classification of the vaginal microbiome. *BMC*

888       *Genomics* **13 Suppl 8**, S17 (2012).

889   93. Cole, J. R. *et al.* The Ribosomal Database Project: improved alignments and new

890       tools for rRNA analysis. *Nucleic Acids Res.* **37**, D141-145 (2009).

891   94. Whitney, J. Testing for differences with the nonparametric Mann-Whitney U test. *J.*

892       *Wound Ostomy Cont. Nurs. Off. Publ. Wound Ostomy Cont. Nurses Soc.* **24**, 12

893       (1997).

894   95. Integrative HMP (iHMP) Research Network Consortium. The Integrative Human

895       Microbiome Project: dynamic analysis of microbiome-host omics profiles during

896       periods of human health and disease. *Cell Host Microbe* **16**, 276–289 (2014).

897   96. Koparde, V., Parikh, H., Bradley, S. & Sheth, N. MEEPTOOLS: a maximum

898       expected error based FASTQ read filtering and trimming toolkit. *Int. J. Comput. Biol.*

899       *Drug Des.* **10**, (2015).

900   97. Li, H. & Durbin, R. Fast and accurate long-read alignment with Burrows-Wheeler

901       transform. *Bioinforma. Oxf. Engl.* **26**, 589–595 (2010).

902   98. Ordoukhanian, P., Nichols, J. & Head, S. R. Primer Extension, Capture, and On-

903       Bead cDNA Ligation: An Efficient RNAseq Library Prep Method for Determining

904       Reverse Transcription Termination Sites. *Methods Mol. Biol. Clifton NJ* **1712**, 253–

905       261 (2018).

906   99. Abubucker, S. *et al.* Metabolic reconstruction for metagenomic data and its

907       application to the human microbiome. *PLoS Comput. Biol.* **8**, e1002358 (2012).

908   100.   Kaminski, J. *et al.* High-Specificity Targeted Functional Profiling in Microbial

909       Communities with ShortBRED. *PLOS Comput. Biol.* **11**, e1004557 (2015).

910

**Table 1: Relative Risk of Vaginal Infections in the Presence of MP1 and MP2**

|  | **Bacterial Vaginosis** | **Trichomoniasis** | **Candidiasis** |
|---|---|---|---|
| *Megasphaera* phylotype 1 | **4.57 (3.76-5.55)** | **0.96 (0.59-1.56)** | **0.52 (0.37-0.74)** |
| *Megasphaera* phylotype 2 | **2.19 (1.79-2.69)** | **4.84 (3.06-7.64)** | **0.54 (0.30-0.96)** |
| *Gardnerella vaginalis* | 6.44 (4.18-9.92) | 2.47 (1.23-4.94) | 0.65 (0.49-0.88) |
| *Prevotella* cluster2 | 5.48 (4.31-6.98) | 2.32 (1.44-3.74) | 0.45 (0.33-0.62) |
| Clostridiales BVAB2 | 4.14 (3.46-4.96) | 1.04 (0.63-1.72) | 0.42 (0.28-0.64) |
| *Sneathia amnii* | 3.97 (3.25-4.84) | 2.93 (1.83-4.70) | 0.48 (0.35-0.68) |
| *Mycoplasma hominis* | 2.57 (2.15-3.08) | 5.80 (3.70-9.09) | 1.07 (0.75-1.53) |
| *"Ca.* Mycoplasma girerdii*"* | 0.88 (0.50-1.55) | 21.00 (13.82-31.91) | 0.71 (0.27-1.87) |

*Relative risk values were calculated based on the customary relative risk formula for MP1, MP2 and taxa known to be associated with BV and trichomoniasis. Relative Risk = (A/A+B) / (C/C+D) where A represents the number of samples where the taxon is present and the participant is diagnosed with the disease, B represents the number of samples where the taxon is present but the participant is not diagnosed with the disease, C represents the number of samples where the taxon is absent but the participant is diagnosed with the disease and D represents the number of samples where the taxon is not present and the participant is not diagnosed with the disease. The relative risk conferred by each taxon is shown with the 95% confidence interval in parentheses for BV, trichomoniasis and yeast infection (Candidiasis). Our total cohort of non-pregnant women (N=3091) was used for this analysis.
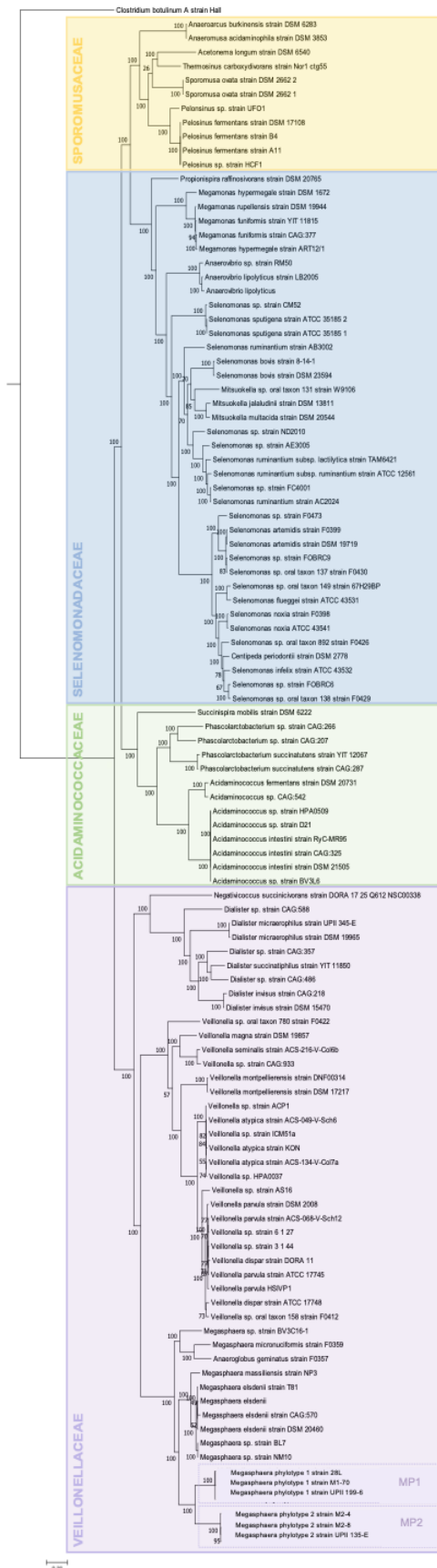
**Figure 1. Maximum Likelihood Phylogenetic Tree of the Class Negativicutes**. A total of 145 orthologous genes from 110 genomes assigned to the class Negativicutes were included in this analysis. *Clostridium botulinum* A strain Hall was designated as the outgroup. This maximum-likelihood phylogenetic tree was generated using 100 bootstrap replicates. Bootstrap values as present at nodes of the tree. Families within the tree highlighted in different colors: Sporomusaceae: yellow, Selenomonadaceae: blue, Acidaminococcaceae: green, Veillonellaceae: purple. MP1 and MP2 genomes are outlined with dotted lines and labeled.

**Figure 2: Maximum Likelihood Phylogenetic Tree of 16S Ribosomal RNA gene**. This maximum likelihood tree was generated using RAxML-HPC with 1,000 bootstraps. Input data were full-length 16S ribosomal RNA gene sequences (nucleotide). Numbers at nodes are indicative of bootstrap support of that node placement. *Dialister micraerophilus* was selected as the outgroup and is a human oral isolate also classified in the family Veillonellaceae. Remaining isolates are colored by their site of isolation: blue- mammalian gut, green- human oral, purple- human vaginal.

44

**Figure 3. Distinctive GC Composition, Codon Preference & Genomic Structure Between Vaginal *Megasphaera* Phylotypes**. a) Differences in both whole genome (blue) and protein-coding (orange) GC composition are shown. b) codon preference is distinct between MP1 and MP2 genomes based on differences in GC composition at specific codon positions (position 1, position 2, position 3 : blue, orange, gray) and in the overall coding GC composition (yellow). c) synteny is conserved within phylotype but lost between MP1 and MP2 genomes. Synteny plots demonstrate structural alignment of genomic content at the amino acid level. Color designates similarity at the amino acid level. The upper panel shows the strong conservation of genomic synteny and protein identity (red color) between two MP1 genomes. The lower panel show massive genome rearrangement and loss of amino acid sequence conservation between a MP1 and a MP2 genome.
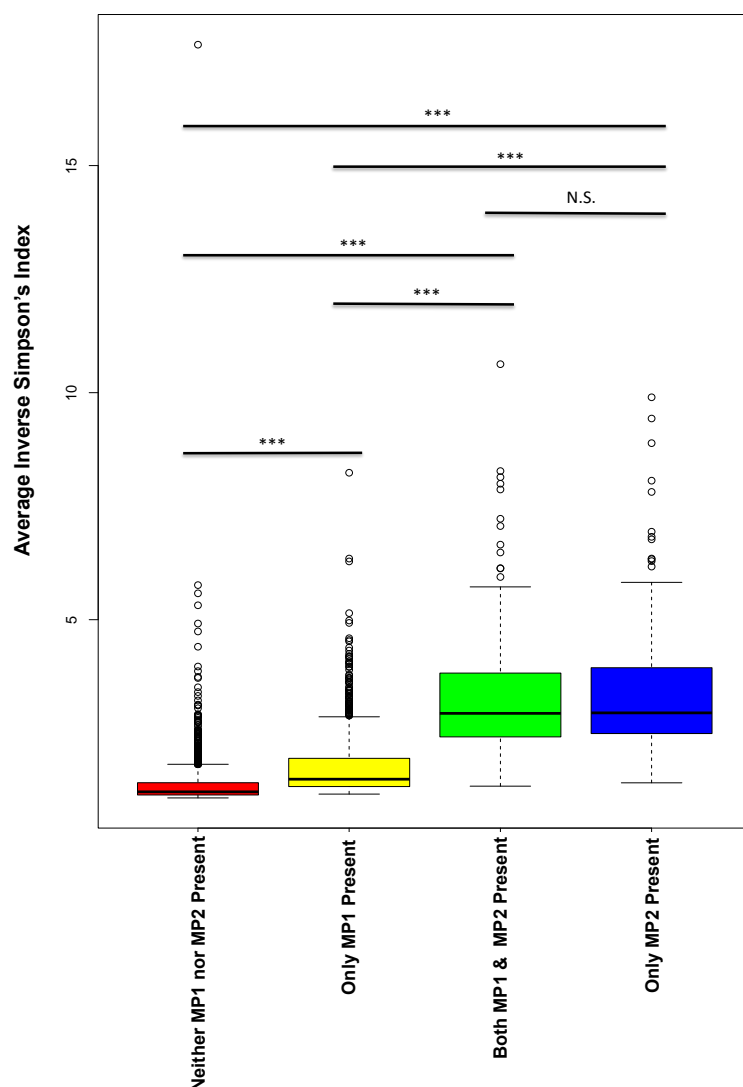
45

**Figure 4. Vaginal *Megasphaera* Phyotypes Associated with Increased Alpha Diversity.** Alpha diversity was measured for vaginal microbiome profiles using the Inverse Simpson's Index, calculated using the 'vegan' package in R. Distribution of Inverse Simpson's Index for each group is shown. Boxes show median and interquartile ranges, with whiskers denoting maximum and minimum values. Outliers are shown as dots. Significance was determined using a two-tailed Student's T-test. Four different groups are shown, samples containing neither MP1 or MP2 (n=1901, red), samples containing MP1 only (n=845, yellow), samples containing both MP1 and MP2 (n=182, green) and samples containing only MP2 (n=163, blue). Taxa were determined to be present in a sample if they comprised greater than or equal to 0.1% of the sample. Samples with MP1 only, MP2 only and both phylotypes all exhibit increased alpha diveristy, with MP2 only samples being the most highly diverse. All comparisons were found to be highly signifcant (p<0.01) with the exception of MP2 only and both phylotypes.
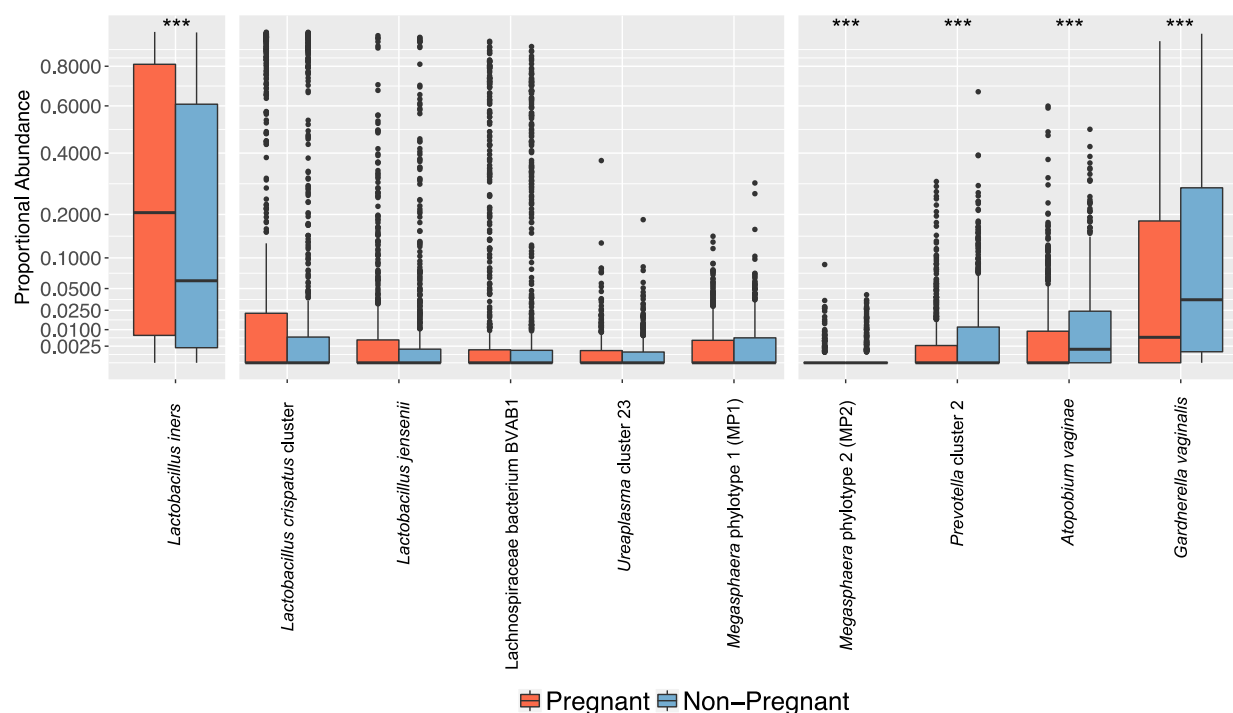
**Figure 5. *Megasphaera* phylotype 1 (MP1) Not Significantly Excluded in Pregnancy**. Results were generated from a cohort of 779 pregnant women case matched 1:1 with non-pregnant controls (N=1558). Using the R packages 'wilcox', a Mann-Whitney U test was performed on all vaginal microbial taxa both present in at least 5% of samples and comprising at least 0.1% relative proportion of the microbiome profile. The R package 'p.adjust' was utilized to correct for multiple testing using the FDR correction. The distribution of proportional abundance across both pregnant (red) and non-pregnant (blue) cohorts are shown. Boxes show median and interquartile ranges, with whiskers denoting maximum and minimum values. Outliers are shown as dots. *Lactobacillus iners* is shown to be significantly more prevalent in the pregnant cohort (q=1.20E-6). *Lactobacillus crispatus* cluster*, Lactobacillus jensenii,* Lachnospiraceae BVAB1, *Megasphaera* phylotype 1 (MP1) and *Ureaplasma* cluster 23 are not significantly different between the two cohorts (q=0.19, 0.23, 0.43, 0.56, 0.26 respectively). *Megasphaera* phylotype 2 (MP2), *Prevotella* cluster 2, *Atopobium vaginae* and *Gardnerella vaginalis* are significantly lower in the pregnant cohort (q=$5.82 \times 10^{-3}$, 6.09E-8, 7.95E-8, 2.82E-7 respectively).
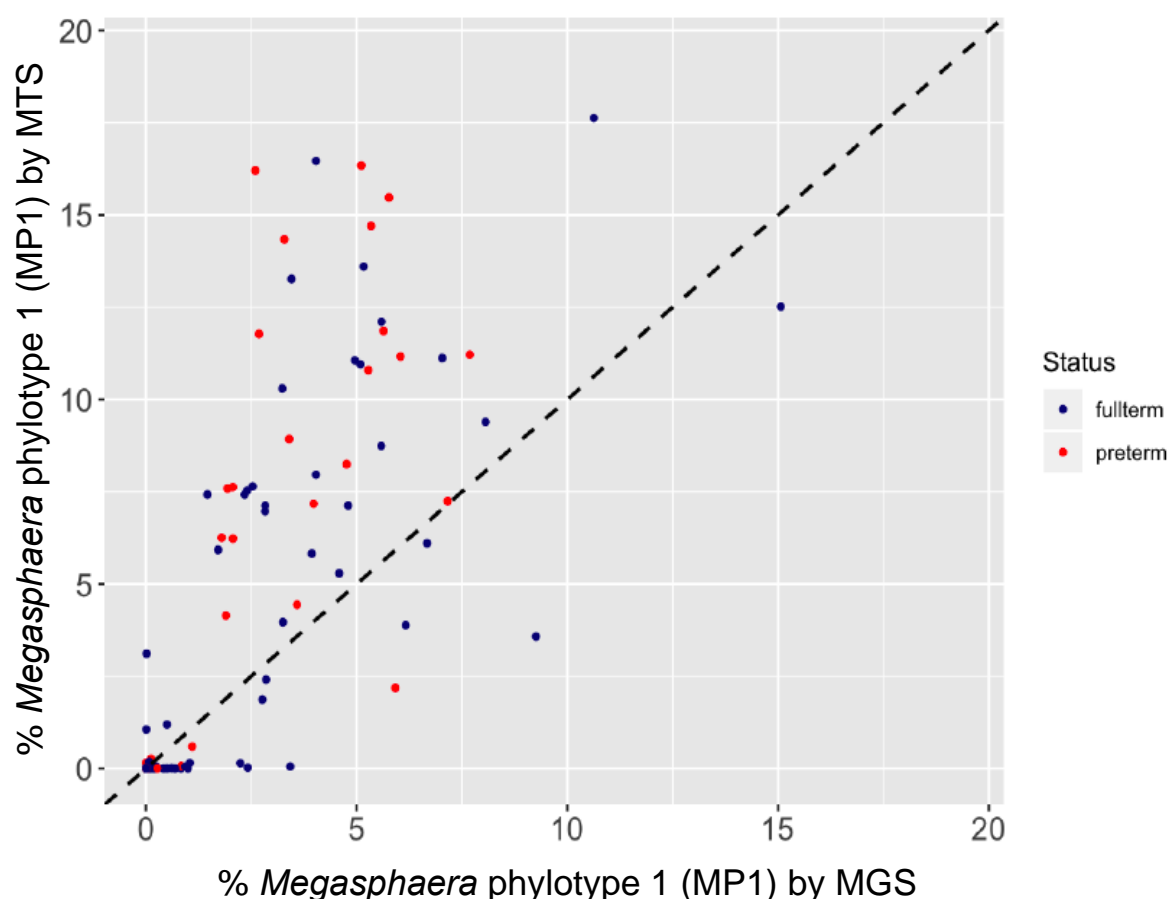
47

**Figure 6. Relationship between *Megasphaera* 16S read abundance and transcript abundance in paired datasets.** Results were generated from samples collected from a cohort of pregnant women that participated in the MOMS-PI study. Samples were processed for whole metagenome microbiomics and transcriptomics. Percent of total transcripts attributed to the taxon of interest is shown on the y-axis. Percent of total whole metagenome sequencing reads attributed to the taxon of interest are shown on the x-axis. Each dot represents an individual sample. The relationship between WMGS and WMTS representation of *Megasphaera* phylotype 1 (MP1) is shown. Figures were generated using the R package 'ggplot'. Data points representing samples from women who went on to deliver full term are shaded blue, while data points representing samples from women who went on to deliver preterm are shaded red. The dotted line extending across the graph diagonally represent the expected 1:1 relationship of WMGS and WMTS- based abundance measures.

**Supplementary Table 1. Genome Characteristics of *Megasphaera* phylotype 1 and *Megasphaera* phylotype 2 Isolates**

|  | M1-70 (MP1) | 28L (MP1) | UPII 199-6 (MP1) | M2-4 (MP2) | M2-8 (MP2) | UPII 135-E (MP2) |
|---|---|---|---|---|---|---|
| **Isolation Location** | VCU | JCVI | JCVI | VCU | VCU | JCVI |
| **Predicted Genome Size (Mb)** | 1.78 | 1.73 | 1.64 | 1.74 | 1.71 | 1.65 |
| **GC Percentage** | 46.33 | 46.05 | 46.37 | 38.94 | 39.09 | 38.88 |
| **Number of Contigs** | 129 | 34 | 45 | 311 | 328 | 49 |
| **N50 length (bp)** | 179993 | 156177 | 100595 | 102411 | 131070 | 64000 |
| **Number of Contigs @ N50** | 4 | 5 | 7 | 6 | 5 | 8 |
| **Transcriptome Size (Mb)** | 1.55 | 1.55 | 1.46 | 1.46 | 1.41 | 1.44 |
| **Transcriptome/Genome Ratio** | 0.87219 | 0.89894 | 0.88832 | 0.83913 | 0.82962 | 0.87491 |
| **Number of Predicted Genes** | 1647 | 1715 | 1457 | 1591 | 1508 | 1510 |

**Supplementary Table 2. Percentage of Conserved Proteins Among Vaginal *Megasphaera* and Closely Related Genomes.** Percentage of Conserved Proteins (POCP) values were calculated based on the method described by Qin et al. A POCP value of less than 50% is indicative that two genomes should be classified to separate bacterial genera. Pairwise POCP values are denoted by color: dark blue- 80-100%, medium blue- 60-80%, light blue- 50-60%, white- less than 50%, likely isolates from distinct genera.

## Supplementary Table 3. Conserved Signature Indel and Conserved Signature Protein Analysis of Vaginal *Megasphaera* Phylotypes

| | *Megasphaera* Phylotype 1 (MP1) strain 28L | *Megasphaera* Phylotype 1 (MP1) strain UPII 199-6 | *Megasphaera* Phylotype 1 (MP1) strain M1-70 | *Megasphaera* Phylotype 2 (MP2) strain UPII 135-E | *Megasphaera* Phylotype 2 (MP2) strain M2-4 | *Megasphaera* Phylotype 2 (MP2) strain M2-8 |
|---|---|---|---|---|---|---|
| **Conserved Signature Indels** | | | | | | |
| Specific to Class *Negativicutes* | | | | | | |
| 3-Isopropylmalate dehydratase, large subunit (1 aa deletion, position 30-71) | Present | Present | Present | Present | Present | Present |
| DNA-directed RNA polymerase, subunit sigma    (1 aa insertion, position 47-73) | Present | Present | Present | Present | Present | Present |
| Specific to Family *Veillonellaceae* | | | | | | |
| GTP disphosphokinase (1 aa deletion, position 441-476) | Present | Present | Present | Present | Present | Present |
| GTP disphosphokinase (1 aa deletion, position 362-403) | Present | Present | Present | Present | Present | Present |
| **Conserved Signature Proteins** | | | | | | |
| Specific to Class *Negativicutes* | | | | | | |
| SELR_02010 | Present | Present | Present | Present | Present | Present |
| SELR_03110 | Present | Present | Present | Present | Present | Present |
| SELR_03270 | Present | Present | Present | Present | Present | Present |
| SELR_05060 | Present | Present | Present | Present | Present | Present |
| SELR_08460 | Present | Present | Present | Present | Present | Present |
| SELR_10260 | Absent | Absent | Absent | Absent | Absent | Absent |
| SELR_10270 | Absent | Absent | Absent | Absent | Absent | Absent |
| SELR_15360 | Absent | Absent | Absent | Absent | Absent | Absent |
| SELR_06480 | Present | Present | Present | Present | Present | Present |
| Specific to Family *Veillonellaceae* | | | | | | |
| MELS_0132 | Present | Present | Present | Present | Present | Present |
| MELS_0206 | Present | Present | Present | Present | Present | Present |
| MELS_0844 | Present | Present | Present | Present | Present | Present |
| MELS_2049 | Present | Present | Present | Present | Present | Present |

*Presence of Conserved Signature Proteins (CSPs) and genomic regions containing Conserved Signature Indels (CSIs) indicative the class Negativicutes and the family Veillonellaceae are shown. All CSIs for both the class and family were detected in MP1 and MP2 genomes. All CSPs indicative of the family Veillonellaceae were also identified. Absent CSPs (3/9 indicative of the class Negativicutes) are denoted in red.

**Supplementary Table 4. 16S Ribosomal RNA Similarity Matrix Among Vaginal *Megasphaera* Phylotypes and Related Genomes.** Pairwise 16S ribosomal RNA similarity was calculated using full length 16S rRNA sequences and the blastn algorithm. Similarity values are denoted by color: red- 98-100%, orange-96-98%, yellow-94-96%,green-92-94%, blue-90-92%. The suggested cutoff for delineating species is 97%.

**Supplementary Table 5. Average Nucleotide Identity Analysis Among Vaginal *Megasphaera* Phylotypes and Related Genomes.** Pairwise Average Nucleotide Identity (ANI) was calculated using a publicly availble script (see Methods). ANI values are denoted by color: yellow- greater than 95%, the suggested cutoff for classifying isolates as the same species, green- 80-94.99%, blue- less than 80% ANI.

**Supplementary Table 6. Predicted Metabolic Differences between MP1, MP2 and Closely Related Bacterial Taxa.** Sheet 1: Genes that distinguish vaginal Veillonellaceae from closely related species are shown including three sections: genes largely conserved in *Megasphaera* and *Anaeroglobus* but lost in MP1 and MP2, genes specific to oral and vaginal strains, and genes specific to MP1 and/or MP2. Sheet 2: Genes distinguishing MP1 and MP2 genomes are shown in three sections: genes present only in MP1, genes present only in MP2, and genes that are variable between the two phylotypes. Genes present in a specific genome are denoted with an 'X'.

**Supplementary Table 7. Vaginal *Megasphaera* Phylotypes exhibit Differential Asscociations with Demographics.** General demographics and clinical measures of the non-pregnant, outpatient cohort (n=3091) are shown. Results are separated into five distinct cohorts: i) the overall cohort (n=3091), ii) participants carrying neither MP1 or MP2 (n=1901), iii) participants carrying MP1 only (n=845), iv) participants carrying MP2 only (n=163) and v) participants carrying both phylotypes (n=182). Counts (left) and percentages (right) are shown for each datapoint with the exception of age and sample pH which are shown as averages.

**Supplementary Table 8. *Megasphaera* Phylotype 1 (MP1) Transcription in Pregnancy.** The Multi- 'Omic Microbiome Study- Pregnancy Initiative (MOMS-PI) cohort was utilized for this analysis. Transcripts were classified using HUMAnN2[98,99] and shortBRED[100] to specific functional pathways. Pathway abundances attributed to MP1 for each sample are shown. Forty-three samples contained MP1 transcripts.

**Supplementary Table 9. *Megasphaera* Phylotype 2 (MP2) Transcription in Pregnancy.** The Multi- 'Omic Microbiome Study- Pregnancy Initiative (MOMS-PI) cohort was utilized for this analysis. Transcripts were classified using HUMAnN2[98,99] and shortBRED[100] to specific functional pathways. Pathway abundances attributed to MP2 are shown. One sample contained MP2 transcripts.

**Supplementary Table 10. Genomes Utilized in Phylogenetic Reconstruction of the Class Negativicutes.** Basic genome statistics acquired from NCBI are shown. The genomes included were utilized in the creation of the class Negativicutes phylogenetic tree. These include genomes classified to the class Negativicutes and available at NCBI

as of January 1, 2015. Genomes were excluded from analysis if they did not contain all 145 orthologous genes needed for the analysis, were low-quality or were significantly different in size or content from other deposited genomes of the same species potentially indicative of a poor asssembly.
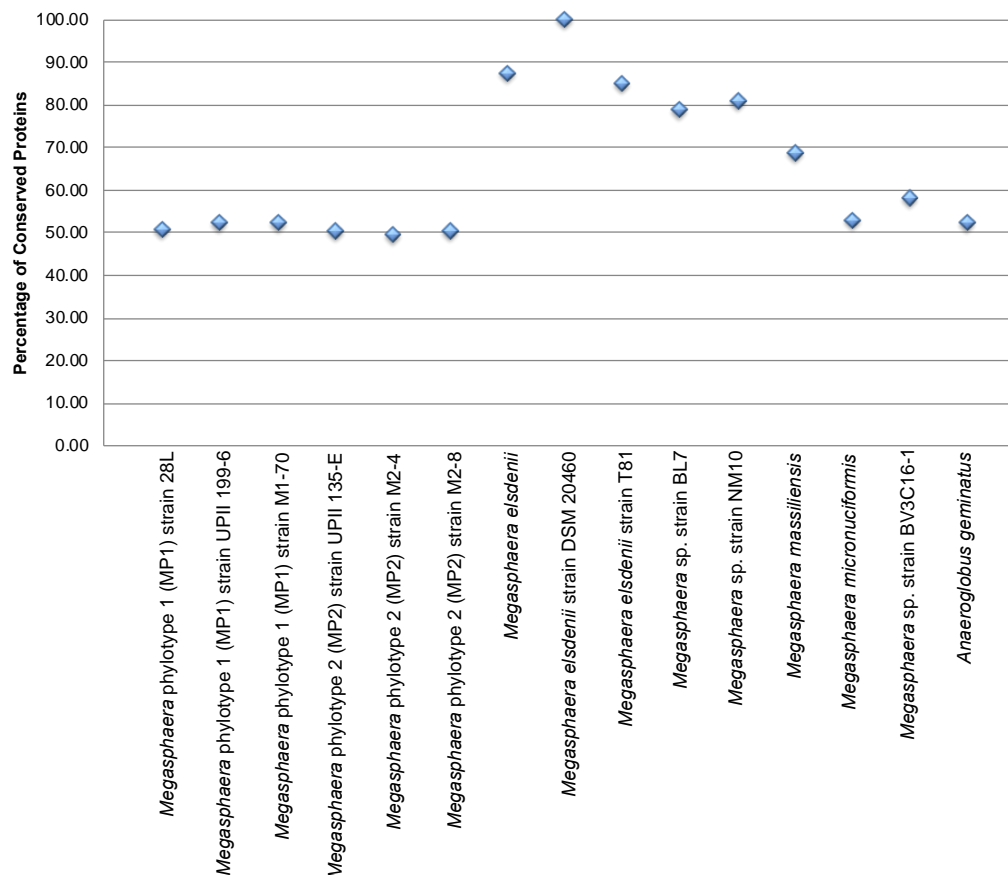
54

**Supplementary Table 11. Supplemented Brain-Heart Infusion Recipe**

| Ingredient | Quantity |
|---|---|
| Brain-Heart Infusion Powder (Oxoid) | 9.25g |
| Yeast Extract | 2.50g |
| Gelatin | 2.50g |
| Dextrose | 0.25g |
| Sucrose | 0.25g |
| Deionized Water | 250mL |

**Supplementary Table 12. Universal 16S rRNA Gene Primers**

| Primer Name | Sequence (5' to 3')[a] |
|---|---|
| 16SF-YM | AGAGTTTGATYMTGGCTCAG |
| 16SF-Bif | AGGGTTCGATTCTGGCTCAG |
| 16SF-Bor | AGAGTTTGATCCTGGCTTAG |
| 16SF-Chl | AGAATTTGATCTTGGCTTAG |
| 1492R | TACCTTGTTACGACTT |

[a] Degenerate bases are underlined. Forward primers were combined in a 4:1:1:1 ratio (16SF-YM : 16SF-Bif : 16SF-Bor : 16SF-Chl).[64,65]

**Supplementary Figure 1. Percentage of Conserved Proteins Analysis versus *Megasphaera* Type Strain.** Pairwise Percentage of Conserved Proteins (POCP) values were calculated based on the methods described in Qin et al., 2014. Shown are the POCP values generated between 15 taxa and the *Megasphaera* type strain *Megasphaera elsdenii* strain DSM20460. POCP values below 50% are the suggested cutoff for delineation of a separate bacterial genus.

**Supplementary Figure 2. Syntenic Comparison of Vaginal *Megasphaera* Phylotypes and the Oral Isolate *M. micronuciformis*.** Full genomes for three MP1, three MP2 and one *Megasphaera micronuciformis* isolate were used for this analysis. Synteny plots demonstrate structural alignment of genomic content at the amino acid level. Color designates similarity at the amino acid level. Synteny is conserved within phylotype as evidenced clear alignment of genomes and protein identity is conserved as well. Between the two phylotypes and in comparsion of *M. micronuciformis*, massive genome rearrangement and loss of amino acid sequence conservation is observed.