# Citizen-science map of the vaginal microbiome

**Sarah Lebeer** ( ✉ sarah.lebeer@uantwerpen.be )

UAntwerpen   https://orcid.org/0000-0002-9400-6918

**Sarah Ahannach**

UAntwerpen

**Stijn Wittouck**

UAntwerpen

**Thies Gehrmann**

UAntwerpen

**Tom Eilers**

UAntwerpen

**Eline Oerlemans**

UAntwerpen

**Sandra Condori**

UAntwerpen

**Jelle Dillen**

University of Antwerp

**Irina Spacova**

UAntwerpen   https://orcid.org/0000-0003-0562-7489

**Leonore Vander Donck**

UAntwerpen

**Caroline Masquiller**

UAntwerpen

**Peter Bron**

University of Antwerp

**Wannes Van Beeck**

UAntwerpen

**Charlotte De Backer**

UAntwerpen

**Gil Donders**

UAntwerpen

**Veronique Verhoeven**

veronique.verhoeven@uantwerpen.be

Biological Sciences - Article

# Citizen-science map of the vaginal microbiome

Sarah Lebeer[1,*,$], Sarah Ahannach[1,$], Stijn Wittouck[1,$], Thies Gehrmann[1,$], Tom Eilers[1], Eline Oerlemans[1], Sandra Condori[1], Jelle Dillen[1], Irina Spacova[1], Leonore Vander Donck[1], Caroline Masquillier[2], Peter A. Bron[1], Wannes Van Beeck[1], Charlotte De Backer[3], Gilbert Donders[4,5,6,°], Veronique Verhoeven[7,°]


*corresponding author

$shared first authors

°shared responsible clinicians


**Affiliations**

[1]Department of Bioscience Engineering, Research Group Environmental Ecology and Applied Microbiology, University of Antwerp, Groenenborgerlaan 171, 2020 Antwerp, Belgium

[2]Department of Sociology, Center for Population, Family and Health, University of Antwerp, Sint-Jacobstraat 2, 2000 Antwerp, Belgium

[3]Department Communication Sciences, University of Antwerp, Sint-Jacobstraat 2, 2000 Antwerp, Belgium

[4]Department of Obstetrics and Gynaecology, University Hospital Antwerp, Drie Eikenstraat 655, 2650 Edegem, Belgium

[5]Regional Hospital Heilig Hart, Kliniekstraat 45, 3300 Tienen, Belgium

[6]Femicare, Clinical Research for Women, Gasthuismolenstraat 33, 3300 Tienen, Belgium

[7]Department of Family Medicine and Population health (FAMPOP), University of Antwerp, Doornstraat 331, 2610 Antwerp, Belgium

**Keywords**

Citizen science / vaginal microbiome / lactobacilli / large-scale remote sampling / population cohort / lifestyle impact

29    **Abstract**

30    The vaginal microbiome is crucial for women's health and reproduction, but its ecology and

31    determinants in the general population are still unclear. This lack of a reference framework

32    hampers much-needed innovations in diagnostics and therapeutics. Here, we remotely

33    mapped the vaginal microbiome of 3,345 women in Western Europe via a citizen-science

34    approach. More than 75% of the vaginal samples were dominated by *Lactobacillus* taxa, but

35    not in discrete community state types. Compositional correlation network analysis validated

36    with public data pointed at six main modules of interacting microbes: a *Lactobacillus*

37    *crispatus-*, *Lactobacillus iners-, Gardnerella-, Prevotella-, Anaerococcus-,* and gut-derived

38    module. In the first module, *Limosilactobacillus* taxa were functionally connected to *L.*

39    *crispatus* and *Lactobacillus jensenii*. This module was positively associated with the luteal

40    phase of the menstrual cycle and negatively with the number of vaginal complaints, while the

41    *Gardnerella*-module was associated with discharge and increasing age. Contraceptives with

42    oestrogen correlated with higher levels of the *L. crispatus-* and less of the *Gardnerella-*

43    module, with the opposite found for a hormonal intrauterine device or having multiple

44    partners. Mothers had lower relative abundance of the *L. crispatus*-module and more

45    *Bifidobacterium, Lactobacillus gasseri* and *Streptococcus*. Other covariates such as BMI,

46    menstrual pads and cups, smoking and dietary habits were also associated with the microbial

47    constellation. These findings suggest that lifestyle interventions have potential to improve

48    vaginal health when combined with dedicated therapies.

## Introduction

The vaginal microbiome plays a central role in women's health and reproduction, but detailed knowledge about its general ecology and the host-side determinants of its composition is lacking. For more than a century, the vagina has been considered a rather simple ecosystem characterized by a low diversity and a high abundance of lactic acid-producing bacteria[1]. In 1892, Döderlein and colleagues described a gram-positive bacterium, as the key bacterium in the vagina[2]. Since then, it has been well established that *Lactobacillus* taxa are the most dominant bacteria in female populations from European and Asian[3–5]. The dominance of these lactobacilli in the vagina is linked to health: when disrupted by an overgrowth of anaerobic bacteria such as *Gardnerella vaginalis* during bacterial vaginosis (BV), or because of inflammation during aerobic vaginitis (AV) or pelvic inflammatory disease (PID)[6,7], an increased susceptibility to conditions such as sexually transmitted diseases[8–10] and adverse reproductive outcomes[11,12] is observed. In 2020, the taxonomy of the family *Lactobacillaceae* was significantly revised[13]. This was an important taxonomic update, as it revealed that the typical vaginal species all belong to the same genus: the *Lactobacillus* genus *strictu sensu*. In addition, the update highlighted the evolutionary distances to other lactobacilli such as *Lacticaseibacillus rhamnosus*, *Lactiplantibacillus plantarum* and *Limosilactobacillus reuteri* that are commonly studied as gut probiotics[13].

With the advent of amplicon sequencing, the vaginal microbiome has been generally described based on five vaginal community state types (CSTs)[3]. *L. crispatus* is dominant in CST I, *L. gasseri* in CST II, *L. iners* in CST III and *L. jensenii* in CST V. CST IV is not dominated by *Lactobacillus*, but rather a mix of more facultative or strict anaerobes such as *Gardnerella, Atopobium, Prevotella,* and *Finegoldia*[14]. This CST IV is found in asymptomatic women but is

72　more associated with dysbiosis and problems such as BV. The recent VALENCIA (VAginaL

73　community state typE Nearest CentroId clAssifier) study proposed thirteen CSTs, based on

74　meta-analysis of 1,976 women from different study cohorts, with particularly extra

75　subdivisions for this CST IV[15]. The CST framework has been very useful to simplify high-

76　dimensional microbial community datasets and facilitate statistical analyses. However, it is

77　currently unclear how well the vaginal CSTs reflect the inherent biology.

78　To better understand the ecology and function of vaginal lactobacilli and other microbiome

79　members and to better design diagnostic and therapeutic options for vaginally associated

80　diseases, more reference datasets are also necessary. So far, female populations in North

81　America, Scandinavia and South-Africa have been mainly characterized[3,15–18], while there

82　seems to be a vast knowledge gap on the vaginal microbiome in other populations. Moreover,

83　other valuable information can come from human-animal comparisons[19]. Humans appear to

84　be the only animals with a vagina mostly dominated by *Lactobacillus* taxa under healthy

85　conditions[20,21]. This unique phenomenon is at present not yet well understood, but the typical

86　hormonal fluctuations throughout the menstrual cycle, particularly estrogen[14]; the glycogen

87　accumulated in the vaginal epithelial cells[22]; the typical human diet since agriculture was

88　introduced[23]; and the strong antimicrobial capacity of lactobacilli that protect the limited

89　offspring of humans from infections[20] have all been suggested to play a role. A detailed

90　mapping of lifestyle and personal characteristics in relation to the vaginal microbiome can aid

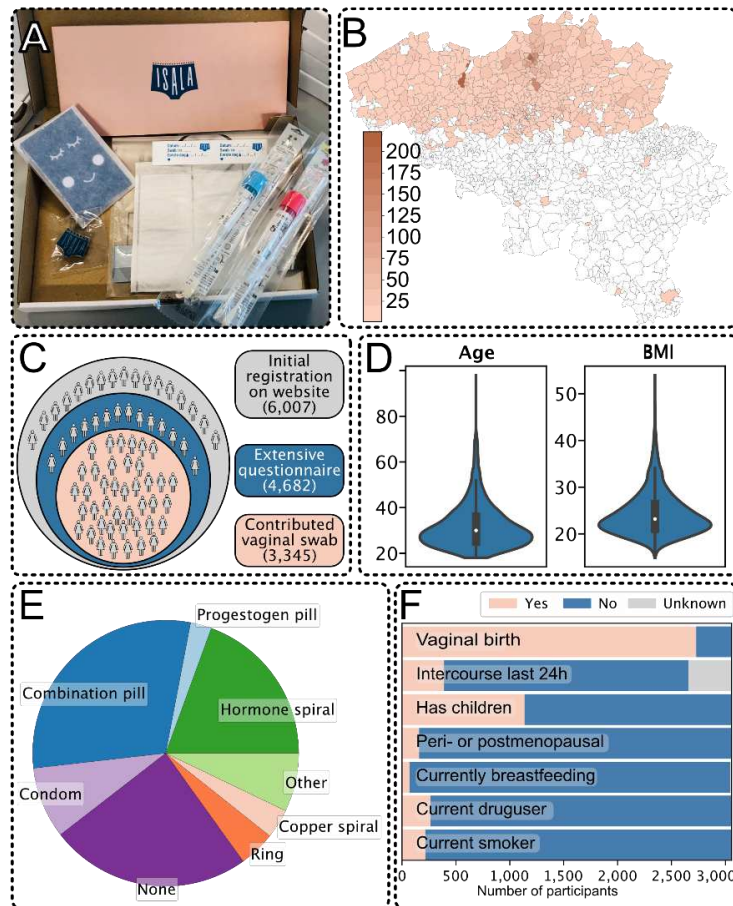91　to better understand the unique build-up of the human vaginal microbiome.

92　In this citizen science-based self-sampling study, we mapped the vaginal microbiome in a

93　large cohort in Belgium, with a particular focus on the prevalence and abundance of key taxa

94　of the lactobacilli, and their association with life-course and lifestyle factors. Two self-

95    collected vaginal swabs were donated by 3,345 women ranging from 18 to 98 years old: one

96    for 16S rRNA amplicon sequencing and one for culturing and metabolic analyses. The project

97    was named 'Isala', after Isala Van Diest (1842-1916), honoring the very first female doctor in

98    Belgium.

99    **Results**

100   **Citizen Science-based study cohort.** The call for participation was launched in Belgium

101   (Western Europe) in March 2020. Within ten days, 6,007 participants registered on the

102   website and registrations were closed (https://isala.be/en/). A total of 4,682 of the original

103   registrants completed the questionnaire with an average completion time of 49 minutes and

104   received the self-sampling kit (Figure 1A-B). The sole exclusion criteria were pregnancy and

105   being younger than 18 years. Of the participants that filled in the questionnaire, 3,345

106   provided two vaginal swabs, allowing microbiome, culturomics and metabolomics analyses

107   (Figure 1C). The mean age and body mass indexes (BMIs) of the included participants were

108   31.8 +- 9.5 years and 24.3 +- 4.6 kg/m$^2$, respectively (Figure 1D).

109   The call was directed towards the general female population outside a clinical setting. Indeed,

110   69.7 % of the women did not report a single vaginal health symptom at the time of sampling

111   based on the questionnaire data (Table S3). 18.3% had one self-reported vaginal symptom,

112   ranging from redness, dryness, odor, increased and/or discoloration of discharge, pain during

113   intercourse, itching, swelling, burning, to urinary infection. Only 7% and 2.6% reported two,

114   or three symptoms respectively. Nevertheless, more than 50% and 70% of the participants

115   answered to have at least once experienced a fungal infection or bladder infection,

116   respectively, which are prevalences in agreement with previous studies[24,25].

**Figure 1 – Characterization of the Isala study cohort and key physiological, behavioral, lifestyle and environmental factors of the participating women.** (A) The self-sampling kit sent to the participants via the national postal service. (B) Geographical overview of the participants that sent in samples for this project, by overlaying their zip codes on a map of the Flemish region and some cities from the Wallonia region of Belgium. Darker colors represent higher numbers of participants with that specific zip code. (C) An overview of the population cohort that registered within ten days after the first announcement, with their different citizen-science roles to the Isala project: minimal involvement by expressing online interest as potential donor via website and answering five questions on age, pregnancy, contraceptive use, country of living until three years and zip code (gray), partial involvement by filling out the extensive questionnaire (blue) and full involvement as donors and 24h follow-up questionnaire (pink). The distribution of a selection of the questionnaire variables: (D) age and BMI, (E) reported contraceptive use of the whole cohort, (F) a subset of the binary variables.

5.2% of the Isala participants were menopausal. 30.2% used a combined oral contraceptive pill, 19.9% a hormonal intrauterine device, 13.1% condoms, 3.7% a copper intrauterine device and 2.5% a progestogen-only pill (Figure 1E). Other forms of contraception (implant, cup, periodic celibacy, sterilization of the participant and/or partner, etc.) were less frequent at 5.7% combined (Table S3). About four out of ten (39.2%) of all women had ever been pregnant. 16.0% reported sexual intercourse within 24h before sampling. 9.1% identified
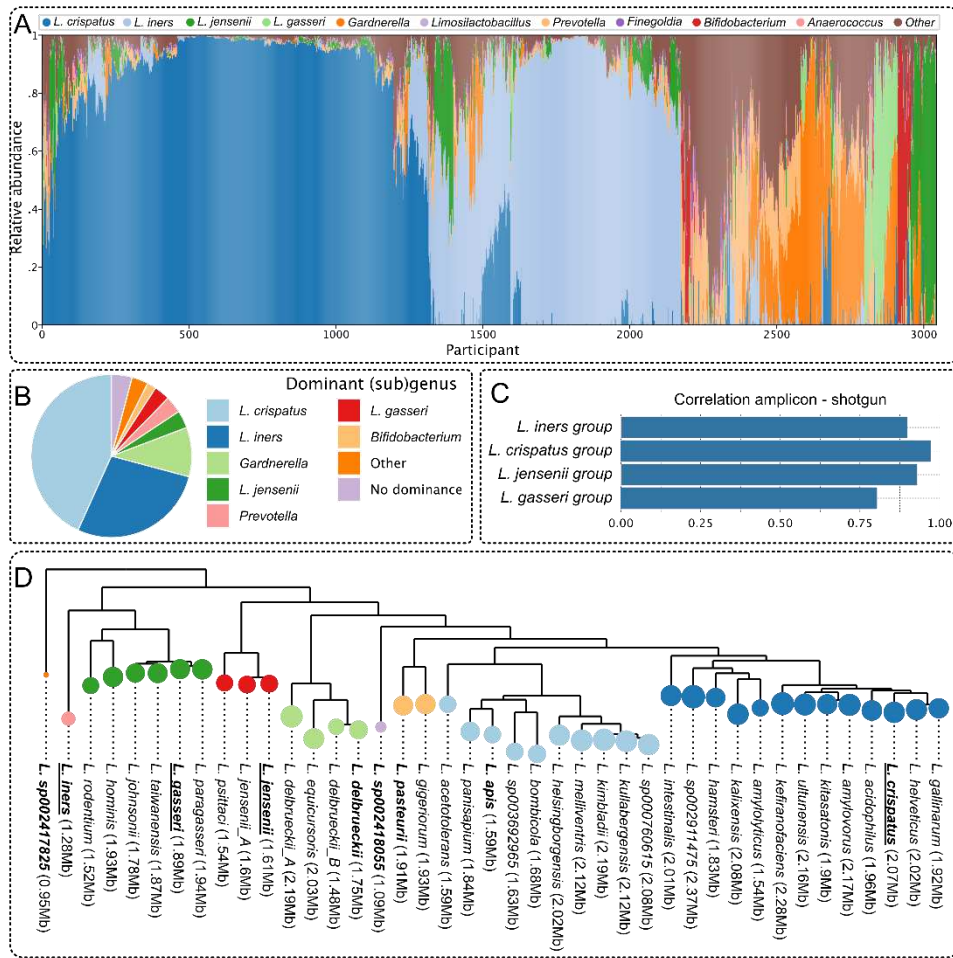
6

136    themselves as a smoker, while 8.6% reported drug use (Figure 1F). As expected, age was

137    significantly correlated with BMI, previous pregnancy, having kids and menopause (Figure S1).

138    4.8% of the participants were not born in Belgium, and 10.0% identified with a culture besides

139    the Belgian one. Ethnicity or race, as previously collected as metadata in US vaginal

140    microbiome studies (Caucasian, African-American, Asian, Hispanic)[26] was not explicitly

141    questioned, since considered not relevant to the Belgian population with its diverse

142    ethnography[27]. 163 participants (5.4%) reported to be part of families below the national

143    poverty threshold, calculated based on the total family income and number of dependents[28].

144    **Dominance of *Lactobacillus* taxa.** 3,345 fully involved Isala donors delivered vaginal samples

145    between July and October 2020, of which 3,196 (96.6%) passed quality control based on

146    estimated DNA concentrations. The high-quality samples totaled over 82 million high-quality

147    V4 16S rRNA read pairs, ranging from 2,126 to 376,242 read pairs per sample with an average

148    of 25,909. Read pairs were merged and denoised into a total of 4,972 unique Amplicon

149    Sequence Variants (ASVs). Short-read 16S rRNA gene sequencing studies generally do not

150    allow for species-level identification[29]. This also applies to many vaginal species: for example,

151    the species *L. jensenii* and *Lactobacillus mulieris* both occur in the vagina, but cannot be

152    discriminated using 16S rRNA gene regions[30]. To be able to analyze the data at the functionally

153    interpretable genus level, but still be able to discriminate between the "big four" vaginal

154    *Lactobacillus* species, the *Lactobacillus* genus was divided into subgenera based on a high-

155    quality core genome phylogeny (Figure 2C-D and Figure S2). This resulted in nine subgenera,

156    four of which are known to be associated with the vagina: the *L. crispatus* group, *L. iners*

157    group, *L. jensenii* group and *L. gasseri* group. To validate this subgenus-level classification

158    approach, shotgun metagenomic sequencing was done for a subset of samples (n = 18, Figure

159    2C). For the four subgenera containing the four typical vaginal *Lactobacillus* species, the

160    relative abundance correlations between the methods were remarkably large (Figure S3).

161    For each sample, the dominant (sub)genus was then determined as the (sub)genus with the

162    largest relatively abundance over 30%. Employing these criteria, the *L. crispatus* group (163

163    ASVs) dominated the largest number of samples (43.2% of the participants), followed by the

164    *L. iners* group (120 ASVs) (27.7%) and *Gardnerella* (49 ASVs) (9.8%). Several smaller dominant

165    taxa also occurred, namely the *L. jensenii* group (54 ASVs) (3.5%)*, Prevotella* (421 ASVs) (3.4%),

166    the *L. gasseri* group (56 ASVs) (3.2%), *Bifidobacterium* (18 ASVs) (1.8%) and *Streptococcus* (52

167    ASVs) (1.2%) (Figure 2A-B).

168    Because of the citizen-science nature of the project, the personal vaginal microbiome profiles

169    were communicated to the participants before the submission of this manuscript (Figure S4

170    and https://isala.be/en/results/). Participants received information about the top eight taxa

171    in the dataset accompanied by information for non-microbiology experts (Figure S5). A

172    feedback questionnaire (n = 2,000) showed that 83% of participants who received their

173    results, perceived them as easy to interpret and 99.6% of participants would volunteer again
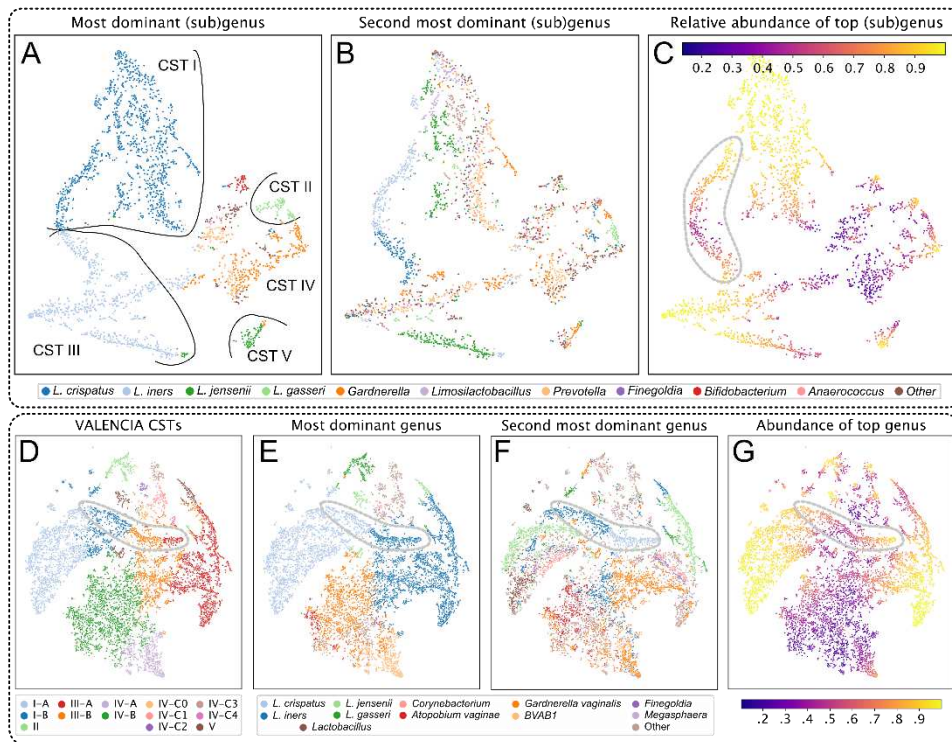
174    in future Isala endeavors.

**Figure 2 - Overview of the most abundant taxa in the vaginal microbiome of the Isala cohort, with particular focus on the *Lactobacillus* taxa.** (A) Stacked bar chart describing the microbiome composition of all participants in the study in terms of the 10 most abundant taxa. (B) Occurrence of the most dominant taxa in the vaginal microbiome of the Isala cohort based on the highest taxonomic resolution possible with our available data. Dominance was defined as the most abundant taxon that constituted at least 30% of the profile. "Other" refers to the number of samples where a different (sub)genus was dominant from the seven that are shown; "no dominance" refers to the number of samples where not a single (sub)genus reached at least 30% abundance. (C) Validation of the 16S amplicon sequencing pipeline, including classification to *Lactobacillus* subgenera, with shotgun sequencing data (n = 18). For the "big four" *Lactobacillus* subgenera, the spearman correlations between their relative abundances in the amplicon and shotgun samples are shown. (D) Maximum-likelihood phylogeny of species of the genus *Lactobacillus* inferred from the amino acid sequences of 100 single-copy core genes. Colors indicate the nine custom-defined subgenera used in this study. Bold tip labels indicate representative species of the subgenera. Species names were taken from the Genome Taxonomy Database[31], which splits species that are very diverse, yielding e.g., *L. delbrueckii_A* and *L. jensenii_A,* the latter recently identified as *L. mulieris*[32]. The size of the circles reflects the genome size of representative genomes of the species (with the average genome size also put between brackets).

195 **Vaginal community structure.** To enable a detailed map of the different constellations of the

196 vaginal microbiota in our cohort, samples were embedded in a two-dimensional t-SNE

197 space[33]. t-SNE projects a high dimensional space into a low dimensional space while aiming

198 to preserve inter-sample distances, placing higher weight on smaller distances to preserve

199 sample neighborhoods. This allows a better global representation of the diversity compared

200 to other commonly used approaches such as PCoA plots[33]. This t-SNE plot was annotated with

201 the two most dominant taxa per sample (Figure 3A-B). Several high-density regions were

202 observed in this two-dimensional representation that broadly corresponded to the five

203 previously described CSTs[3], but these high-density regions were connected by intermediate

204 regions (Figure 3A). A clear example was provided by the *L. crispatus* and *L. iners* high-density

205 regions, which were connected by samples with *L. crispatus* and *L. iners* as the two most

206 abundant taxa. This was the case for 454 samples, of which 22% contained *L. crispatus* and *L.

207 iners* in near-equal proportions (Figure 3C, gray dashed enclosure). This observation suggests

208 that the previously described CSTs are not distinct possibilities in vaginal community

209 composition. This is especially apparent when visualizing the samples based on the second

210 most dominant (sub)genus (Figure 3B) and the relative abundance of the top (sub)genus

211 (Figure 3C). Intermediate regions can be observed in which at least two subgenera are co-

212 dominant, with the same patterns observed in the datasets aggregated in the VALENCIA

213 study[15] (Figure 3D-G). As in the Isala data, samples dominated by *L. iners* and *L. crispatus* at
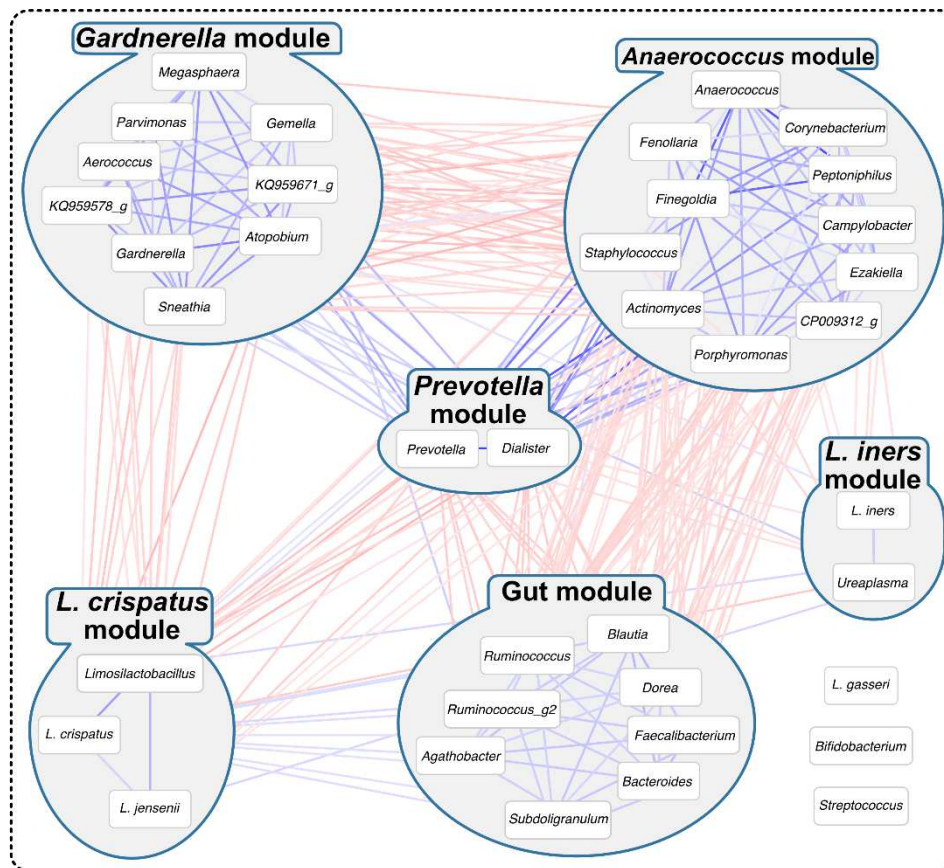
214 near equal abundances were also observed here.

215

**Figure 3 – Vaginal microbiome structure of the Isala cohort.** (A) t-SNE plot of microbiome samples in the Isala study. Embedding colored by the most abundant (sub)genus. Broad community state types (CSTs) are delineated with black lines, except CST IV, which is composed of the remaining samples. (B) Samples are colored by the second-most abundant (sub)genus. (C) Samples are colored by the largest relative abundance level in each sample. (D) Structure of the vaginal microbiome of the VALENCIA public dataset. A t-SNE plot of all microbiome samples of the VALENCIA dataset (multi-temporal samples per participant included), colored by the 13 CSTs presented in that paper. CST I—*L. crispatus* dominated (A high relative abundance, B lower relative abundance), CST II—*L. gasseri* dominated, CST III—*L. iners* dominated (A high relative abundance, B lower relative abundance), and CST V—*L. jensenii* dominated. CST IV-A - *Candidatus Lachnocurva vaginae* (BVAB1) with some *G. vaginalis.* CST IV-B - *G. vaginalis* with low relative abundance of *Ca. L. vaginae*. CST IV-C0 - *Prevotella*, CST IV-C1—*Streptococcus,* CST IV-C2—*Enterococcus* dominated, CST IV-C3—*Bifidobacterium* dominated, and CST IV-C4—*Staphylococcus* dominated. Samples of the VALENCIA dataset colored by (E) the most dominant genus, (F) the second most dominant genus, (G) and by the largest relative abundance level in each sample. The branching point between *L. crispatus* dominated and *L. iners* dominated samples is indicated with a grey line. Of note, BVAB1 corresponds to genus EU728721_g in the Isala dataset, where it only occurred in 1.4% of the participants (not visualized in panel A-B because not in top 10).

The correlation between taxa abundances was investigated with SparCC, considering the compositionality of the relative abundance data[34]. Six main modules of intercorrelated taxa were determined (Figure 4). The first module contained the *L. crispatus group*, *L. jensenii group,* and *Limosilactobacillus*. Correlations between the taxa in this module were weakly positive (r = 0.18 – 0.40). A second module was assigned to a group of taxa that included

239    *Gardnerella, Sneathia*, *Atopobium* and *Aerococcus* (*Gardnerella* module, r= 0.11-0.5). A third

240    module contained the relatively strongly correlated *Anaerococcus*, *Peptoniphilus* and

241    *Finegoldia* taxa (*Anaerococcus* module, r= 0.1-0.71), together with some more weakly

242    correlated taxa such as *Staphylococcus.* A fourth module was composed of *Prevotella* and

243    *Dialister* (r=0.78), which jointly correlated positively with both the *Gardnerella* and

244    *Anaerococcus* modules, while the latter two were negatively correlated with each other. A

245    fifth module was composed of taxa associated with the gut, including *Ruminococcus,*

246    *Bacteroides,* and *Subdoligranulum* (Gut module, r=0.16-0.28). Interestingly, the Gut module

247    was positively correlated with the *L. crispatus* module. Finally, the sixth main module

248    constituted the *L. iners* group and the genus *Ureoplasma*. A few taxa did not show any strong

249    correlations with other taxa, notably *Bifidobacterium, Streptococcus* and the *L. gasseri* group.

250    Yet, when computing SparCC correlations in the VALENCIA dataset, we identified a striking

251    concordance with the modules identified in the Isala dataset (Figures S6 and S7).  In both

252    datasets, the *L. crispatus* module showed moderately negative correlations to the taxa in the

253    *Gardnerella*, *Anaerococcus* and *Prevotella* modules (-0.22, -0.15, and -0.27 respectively),

254    which is in line with the previously documented inhibitory capacity of *L. crispatus*-dominated

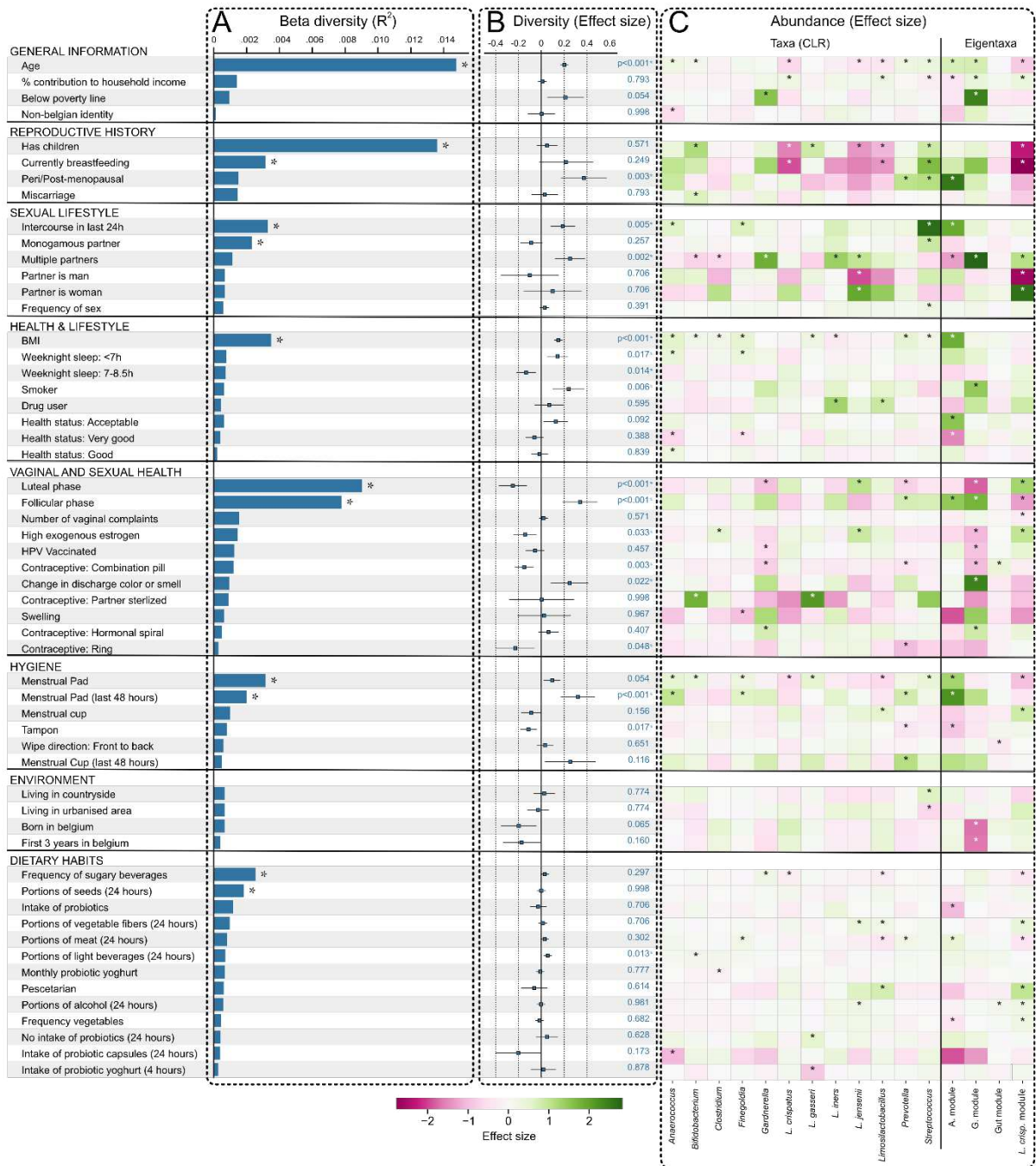255    communities against these potential vaginal pathobionts [35,36].

256

**Figure 4 – Six main modules of interacting microbes as defined by a compositional correlation analysis.** Modules are enclosed in gray. Positive correlations in blue, negative correlations in red. Thickness of the line indicates the strength of the correlation. Exact correlations are given in Figures S6 and S7.

261  Our analysis also pointed at a strong correlation between the genus *Limosilactobacillus* and

262  both the *L. crispatus* and *L. jensenii* groups (which were also positively correlated with each

263  other). *Limosilactobacillus* taxa did not show a high average relative abundance (0.4%) in our

264  dataset, but had a surprisingly high prevalence of 47.8% (Figure 2A and Table S1). Based on a

265  case-by-case ASV sequence comparison with a 16S reference database, we could assign the

266  ASVs classified as *Limosilactobacillus* to one of three groups within the genus: a *Lactobacillus*

267  *reuteri* group, the species *Limosilactobacillus coleohominis* and the species

268  *Limosilactobacillus fermentum*. The *L. reuteri* group contained the species *Limosilactobacillus*

269  *reuteri*, *Limosilactobacillus vaginalis* and five other species that are not known to occur in the

270  human vagina. We found a prevalence of 43.7% for the *L. reuteri* group, 11.5% for *L.*

271  *coleohominis* and 4.1% for *L. fermentum*. In addition to our large dataset with amplicon-

272  sequenced samples, we also inspected the 264 vaginal metagenomes of the VIRGO metastudy

273  for the presence of *Limosilactobacillus* species. The most prevalent species were *L.*

274  *coleohominis* (25%), *L. vaginalis* (20%) and *L. fermentum* (1%)[37]. *L. fermentum* was most

275  frequently cultured from a subset of 592 vaginal swabs, with even more isolates obtained

276  than for *L. crispatus* and *L. jensenii* based on standard growth conditions for lactobacilli (Table

277  S1). Overall, culture of the vaginal lactobacilli was cumbersome under the standard conditions

278  and remains to be further optimized.

279  **Impact of host covariates on the vaginal microbiome.** We then analyzed the association of

280  personal data with key features of the vaginal microbiome (Figure 5). As an alternative to

281  reducing the dimensionality of the microbial community data through a classification into

282  CSTs, alpha and beta-diversity metrices, twelve individual (sub)genera of interest and

283  eigentaxa (see Methods) of the four largest modules of intercorrelated taxa were selected for

284  association testing. The functional relevance of this latter approach was confirmed by the

285  association observed between change in discharge and an increase in the *Gardnerella*-

286  module, but not with specific taxa. Similarly, a lower relative abundance level of the *L.*

287  *crispatus*-module was associated with an increased number of vaginal complaints specifically.

288  Considering age had the largest effects, the data were also adjusted for this parameter.

289



290
**Figure 5 - Statistical analysis of the association of different personal, reproductive, lifestyle, health, hygiene, environmental and dietary factors with the vaginal microbiome space.** Each panel displays effects on different levels of the microbiome: (A) the effect on the beta-diversity between the samples (Adonis test), (B) the effect on the alpha-diversity of the samples, (C) the effect on the abundances of specific taxa and on the eigentaxa of the modules discovered in the SparCC correlation analysis. The A and G modules refer to the *Anaerococcus* and *Gardnerella* modules, respectively. Asterisks represent significant associations (FDR adjusted and using a threshold of 0.05; white and black asterisks are merely for visualisation purposes). The number of samples for each question was almost the entire study (n = 3,043 participants). Due to missing data or specific comparisons, this can deviate, and detailed counts are provided in Table S3.

301 Besides age, having had children had the strongest association with beta-diversity, explaining

302 1.4% of the microbiome variation. It was significantly negatively associated with the

303 abundance levels of *L. crispatus*, *L. jensenii* and *Limosilactobacillus* (the *L. crispatus*-module),

304 and positively with *Bifidobacterium, L. gasseri,* and *Streptococcus*. Breastfeeding at the time

305 of sampling was correlated with beta-diversity, lower relative abundance of *L. crispatus* and

306 *Limosilactobacillus* and higher levels of *Streptococcus*. Being "peri- or post-menopausal" did

307 not show a significant association with the beta-diversity, but it was correlated with an

308 increased alpha-diversity and levels of *Streptococcus, Prevotella* and the *Anaerococcus*-

309 module. Having had intercourse in the last 24 hours was associated with a higher alpha

310 diversity, and higher levels of *Anaerococcus, Finegoldia,* and in particular *Streptococcus.* We

311 also investigated the associations of partnership with the vaginal microbiome. Compared to

312 not being sexually active, having a monogamous relationship correlated with the beta-

313 diversity and higher levels of *Streptococcus*, but no associations were noted for the alpha-

314 diversity. However, having multiple partners was linked with a higher alpha-diversity and

315 higher levels of the *Gardnerella*-module, but also higher levels of the *L. crispatus*-module, and

316 less of the *Anaerococcus*-module. Having a male partner was associated with lower levels of

317 *L. jensenii* and the *L. crispatus*-module, compared to having a female partner. The impact of

318 the stage of the menstrual cycle was evaluated for pre-menopausal participants not taking

319 any related hormonal contraceptives, with the follicular phase starting on the first day of

320 menstruation and the luteal phase after ovulation (Figure S8). As expected, the follicular

321 phase was associated with higher alpha-diversity, together with lower levels of the *L.*

322 *crispatus*-module and higher levels of *Prevotella* and the *Gardnerella-* and *Anaerococcus*-

323 modules, compared to the ovulation and luteal phase. The opposite was true for the luteal

324 phase (compared to the ovulation and follicular phase). Combining the data for

325   contraceptives with a high predicted exogenous estrogen level (combination pill, vaginal ring

326   or patch) showed an association with an increase in the *L. crispatus*- module and less of the

327   *Gardnerella*-module. The oral combination contraceptive pill, which disrupts the natural cycle

328   and contains estrogen and progestin[23], correlated with lower alpha-diversity, lower relative

329   abundances of *Prevotella* and *Gardnerella* but higher levels of the gut taxa module. Use of a

330   ring contraceptive was linked to a significantly lower alpha-diversity and lower levels of

331   *Prevotella*. Use of a hormonal intra-uterine device (containing only progestin) was associated

332   with more of the *Gardnerella*-module. Having been vaccinated against HPV was linked to

333   lower levels of the *Gardnerella*-module. Furthermore, we also observed associations for

334   menstrual hygienic products, with a menstrual cup appearing more beneficial for the *L.*

335   *crispatus*-module and pads being more associated with an increased alpha diversity. The

336   menstrual pads also significantly reduced the *L. crispatus*-module and increased the

337   *Anaerococcus*-module, especially when used in the last 48h. Wiping the vulva from front to

338   back after a bathroom visit was associated with lower levels of the gut taxa module in the

339   vagina.

340   Among the general health and lifestyle factors that were questioned, the largest effect was

341   BMI, which was significantly associated with the beta-diversity, higher alpha-diversity, and

342   higher levels of bacteria in the *Anaeroccocus*-module. Specific dietary components were also

343   linked with the overall composition and diversity of the vaginal microbiome when adjusting

344   for age. The consumption of sugary beverages was noticeably associated with beta-diversity,

345   and with lower levels of the *L. crispatus* module, while the consumption of light beverages

346   (marketed as diet, sugar-free, zero-calorie or low-calorie) in the last 24h was associated with

347   a significantly higher alpha-diversity and higher levels of *Bifidobacterium*. A high portion of

348   seed consumption was significantly associated with beta-diversity, but not with the specific

349 taxa or modules that we examined. High frequency of vegetable consumption and its

350 associated fibers, particularly in the last 24h, and being pescatarian were associated with a

351 minor increase of *L. crispatus*-module. Ethanol consumption in the past 24h was associated

352 with higher levels of the *L. crispatus*- and gut taxa module. Meat consumption was linked to

353 lower levels of the *L. crispatus*-module, and higher levels of *Prevotella* and the *Anaerococcus*-

354 module. Significantly lower levels of the *Anaerococcus* module taxa occurred when probiotic

355 capsules were consumed in the last 24 hours. In contrast, consumption of probiotic yoghurts

356 in the last 24h was associated with lower relative abundance of *L. gasseri*.

357 Additional lifestyle factors other than diet were also evaluated. Sleeping less than seven hours

358 per weeknight corresponded to a significantly higher alpha-diversity and higher levels of

359 *Anaerococcus* and *Finegoldia*, while sleeping between 7 and 8.5 hours corresponded to a

360 lower alpha-diversity. In addition, smoking was associated with higher alpha-diversity, and

361 higher levels of the *Gardnerella*-module. While taking drugs was not linked to the diversity of

362 vaginal samples, it was linked to higher levels of *L. iners* and *Limosilactobacillus* . Income

363 inequality within couples did not show a significant effect on the vaginal microbiome but

364 being below the Belgian poverty threshold was linked to a higher alpha-diversity, and in

365 particular, higher levels of *Gardnerella.* Being born in Belgium and living there for the first 3

366 years was associated with significantly lower levels of the *Gardnerella-module.* Furthermore,

367 living in a more urbanized/polluted area (i.e., city center, village center, busy road, industrial

368 zone) versus suburban/countryside environment (i.e., residential area, rural area, green

369 zone/recreation zone) was associated with lower versus higher levels of *Streptococcus*.

370      All significant factors mentioned above could explain 8.01% of the variation in the vaginal

371      microbiome, compared to 7.63% of the variations explained by covariates in a related study

372      on the gut microbiome in the Belgian population[38].

373      **Discussion**

374      The Isala citizen science project on the vaginal microbiome was inspired by a strong need for

375      a better understanding of the vaginal microbiome outside a clinical setting. The enthusiasm

376      of participants willing to donate intimate samples is in line with the current trend of more

377      women taking their health into their own hands. The fact that our study was fully remote had

378      both advantages and limitations. No blood samples, clinical exams or host genetics data could

379      be obtained, but the fully remote setting and large online questionnaire also provided us with

380      unique opportunities to gain widespread access to samples and intimate data. Other inherent

381      limitations of our study cohort were the slight bias towards a high socioeconomic status, like

382      many other citizen science studies[39,40], and the fact that we had to rely on only one timepoint

383      sampled per participant. On the other hand, the fact that intimate self-sampling could be

384      done in the privacy of the home setting had a positive impact on the number of women willing

385      to participate, resulting in a large, diverse set of samples with sufficient variation to study key

386      parameters such as age, BMI, menstrual cycle, contraceptive use, menopausal status,

387      obstetrical parameters, sexual and vaginal health, diet, income, and sleeping habits. The fact

388      that the analysis of all samples was done within the same lab and a small timeframe

389      minimized the technical variability. Taken together, this study set-up enabled us to obtain

390      novel insights in the average vaginal microbiome constellation of this self-reported healthy

391      Western European population.

392   The first key finding of this work was the high number of participants with a dominance of

393   *Lactobacillus* in this Western-European population cohort: 75% of the women were

394   dominated by *Lactobacillus* taxa, in particular by taxa belong to the *L. crispatus* and *L. iners*

395   group, comparable to similar studies[3,20]. Subgenus or group level classification was preferred

396   to better reflect the diversity in ASVs than generally reported. The *L. crispatus* group (163

397   ASVs) was detected in 43.2% of the participants. *L. iners* was dominant in 27.7% of the

398   participants. As we and others have previously reviewed, *L. iners* has an ambiguous role in

399   the vagina[41]. The fact that we found *L. iners* to be so prevalent in complaint-free women

400   suggests that it is often probably rather a friend than a foe in healthy women. Yet, we

401   observed a high diversity of ASVs for *L. iners* (120 unique ASVs), in line with previous

402   suggestions of different clones of *L. iners* with distinct functional properties[42]. Similarly,

403   *Gardnerella* was dominant in 9.8% of the Isala women, although it is often considered a

404   pathobiont in the vagina. Yet, the association of *Gardnerella* with symptoms and disease

405   appears to depend on the specific species and strains [43,44], the other members in the vaginal

406   community[45] and the host[46]. This context- and taxon-dependent role of the vaginal bacteria

407   highlights that it is important to capture the diversity of the vaginal ecosystem in the most

408   biologically relevant way. From five[3] to thirteen CSTs[15] have been previously proposed. CSTs

409   often confuse clinicians and researchers, as they have been mainly proposed for statistical

410   and epidemiological purpose[15], and should not be interpreted as stable community state

411   types. With t-SNE embedding analyses, we clearly showed that the vaginal microbiome space

412   is a continuum, highlighting that CSTs should not be interpreted as the existence of fully

413   discrete states of the vaginal microbiome, as is now also increasingly recognized[15,47,48]. For

414   example, the two most abundant taxa, the *L. crispatus* and *L. iners* groups frequently co-

415   occurred in varying  and even equal proportions. As an alternative approach to maximally

416     capture the diversity of the microbial space while still enabling the analysis of associations

417     with as many metadata as possible, we introduced modules of taxa of interacting vaginal

418     bacteria (with positive correlations within and mostly negative correlations between

419     modules), for which we made eigentaxa for correlation analyses. The taxa-taxa correlations

420     likely reflect relevant biologic phenomena including positive or negative microbial

421     dependencies such as cross-feeding[49–51], inhibition via antimicrobial production[52] but also

422     different immune or inflammation states of the host, where different "states" of the host

423     enrich or restrict different bacteria[45]. The fact that we could validate the existence of these

424     modules in another large independent dataset (VALENCIA) highlights their biological

425     relevance and existence independent of our dataset, in contrast to CSTs obtained by

426     hierarchical clustering which are more dataset dependent.

427     The *L. crispatus*-module probably reflects the most common healthy homeostatic state, based

428     on the known associations of these lactobacilli with vaginal health[53] and our own observations

429     of a reduced abundance of this module with increased number of vaginal complaints versus

430     its increase with increasing estrogen levels. Notably, the association between this module and

431     vaginal complaints was lost with the individual taxa, showing the added value of

432     implementing these modules. Another unprecedent finding for this module is the prevalence

433     and possibly stabilizing capacity of *Limosilactobacillus*. This genus was shown to be highly

434     prevalent, with occurrence in almost 50% of the women sampled, and showed to be easier to

435     culture than the classic big four (i.e., *L. crispatus, L. iners, L. gasseri* and *L. jensenii*). Positive

436     interactions between different taxa of lactic acid bacteria are very common in food

437     fermentations where lactic acid bacteria dominate. In yoghurt, for instance, *Streptococcus*

438     *thermophilus* and *Lactobacillus delbrueckii* subsp. *bulgaricus* exchange crucial metabolites, a

439     process called protocooperation[49]. In kefir, it was recently shown that *Lactobacillus*

440    *kefiranofaciens*, which dominates the kefir community, uses kefir grains to bind together all

441    other microbes that it needs to survive[50]. Such mutualistic interactions have also been

442    observed for related *Lactobacillus* taxa within vertebrate hosts. For example, in the rodent

443    gastrointestinal tract, *Lactobacillus johnsonii* needs *L. reuteri* for biofilm formation[54]. It

444    appears plausible that a similar interaction occurs in the vagina between species of the same

445    two genera, where one or more *Limosilactobacillus* species support *L. crispatus* and *L. jensenii*

446    as keystone taxa. Of note, one of the most widely used vaginal probiotics, *L. reuteri* RC-14,

447    has been shown to have the capacity to prevent BV in women with HIV[55,56] and improve the

448    BV cure rate with single dose of tinidazole[57]. Yet, in these previous studies, it is difficult to

449    differentiate the effect of *L. reuteri* RC-14 from the other applied probiotic strain

450    *Lacticaseibacillus rhamnosus* GR-1[57].

451    While the *L. crispatus* module contains presumed health-associated taxa, three of the

452    modules contain taxa previously associated with dysbiosis: the *Gardnerella*-module consists

453    mostly of taxa associated with BV[58,59], while the *Anaerococcus*- and *Prevotella* modules also

454    contain taxa previously associated with BV[45,60], but also with more inflammatory host states

455    such as AV[6,7,61], endometriosis[62] and PID[63]. The negative correlation between the *Gardnerella*

456    and *Anaerococcus* modules is in line with the view that BV and other inflammatory states such

457    as AV are different forms of dysbiosis with different underlying causes[7]. In this light, the

458    positive correlation of the *Prevotella* module with both modules is harder to explain and

459    requires further investigation. Interestingly, the number of different vaginal complaints

460    reported by the participants was not significantly associated with any of the three modules

461    containing taxa known to be dysbiosis-associated, but only with a reduction of *L. crispatus*

462    module taxa. This suggests that the presence of these modules in itself is not sufficient for a

463    dysbiotic state to develop; such a development would require an extra host-side factor such

464    as a lack of immune control (such as sometimes thought for BV[45,64] or the development of an

465    inflammatory state (such as observed in AV)[6,18]. For change in discharge, it is noticeable that

466    we found a clear association with the *Gardnerella*-module, but not with the individual taxa,

467    highlighting again the relevance of microbe-microbe interactions. Similarly, we interpret our

468    observation of a gut taxa module by the existence of a gut-vagina axis, which is not only a

469    source of potential urogenital pathogens but also of beneficial colonizers. For the latter, the

470    positive correlation with the *L. crispatus* module is of particular interest.

471    Having established this update picture of the vaginal microbiome constellation and collecting

472    a large dataset of personal data via questionnaires, allowed us to then perform an in-depth

473    analysis of covariates. We could confirm previously found associations such as for BMI[65], the

474    contraceptive pill[66] and smoking[67]. The fact that in our dataset especially estrogen-containing

475    contraceptives had a positive association with the levels of the *L. crispatus*-module, and were

476    also linked to less of the *Gardnerella*-module, is in a way reassuring, given the fact that it is so

477    widely administered in Western Europe and completely abolishes the spontaneous menstrual

478    cycle. A disruption of the vaginal microbiome does not seem a major side effect of the

479    combination pill, although we and many Isala participants acknowledge the existence of other

480    side effects, including impact on mood and libido[68–70] and increased risk for venous

481    thromboembolism[71,72] , which are important to consider when choosing the personally most

482    suitable contraceptive method. Notably, the association of a progestin-containing IUD and

483    increased *Gardnerella*-module found here could be included in information provided to

484    women choosing this contraceptive method. Our data are in line with clinical data that

485    insertion of a hormonal IUD temporarily increases BV and over time increases *Candida* spp.

486    colonization in the vagina[73], while systemic progestin-only contraceptives appear to have

487    mixed effects on the vaginal microbiome[74].

488   The life event with the most significant impact on the vaginal microbiome was having children

489   or having been pregnant, which correlated with an overall reduction in *L. crispatus, L. jensenii*

490   and *Limosilactobacillus* levels and an increase in *Streptococcus, Bifidobacterium* and *L. gasseri*

491   levels. A higher taxonomic resolution was not possible, but these three genera contain taxa

492   beneficial to babies as initial colonizers of the oral cavity and gut of newborns[75]. It has been

493   previously shown that most women experience a postdelivery disturbance in their vaginal

494   microbiome, characterized by a decrease in *Lactobacillus* species and increase in diverse

495   anaerobes which persisted for up to one year[76]. In our Isala dataset, it was surprising that we

496   observed the signature of reduction in the *L. crispatus*-module and increase in *Streptococcus,*

497   *Bifidobacterium* and *L. gasseri* in all women having biological children, independent of their

498   age.  This suggests that the impact of pregnancy could be long-lasting. We have at present no

499   explanation for this phenomenon, although we do acknowledge we have a rather young

500   cohort (average age 31.8 +- 9.5 years). Of note, breastfeeding women (who recently

501   delivered) showed similar and even stronger associations for reduction in *L. crispatus* and

502   increase in *Streptococcus*. Hormonal and associated sugar-level changes during pregnancy

503   (including lower estrogens during breastfeeding), as well as the cervix shortening could all be

504   involved and provide interesting aspects for further research. Moreover, the fact whether

505   childbirth has taken place by vaginal or abdominal mode (C-section), the latter with or without

506   preceding labor (i.e., secondary or primary C-section), may have played a major role, and

507   remains to be elucidated in further studies.

508   Another intriguing finding of our Isala citizen-science study is how dietary choices could have

509   a small, but significant impact. For example, intake of vegetable fibers, alcohol consumption

510   and being a pescatarian had a  significant beneficial impact on the *L. crispatus*-module, while

511   drinking sugary beverages had a negative impact. These associations should obviously be

512    interpreted with care and not taken as one-on-one directions towards lifestyle

513    improvements. Alcohol consumption, for example, was associated with a higher abundance

514    of the *L. crispatus* module, but has an established detrimental impact on the gut

515    microbiome[77]. By contrast, limiting intake of sugary drinks appears a lifestyle intervention

516    that benefits multiple habitats that make up the human body. Another intriguing finding was

517    the different associations found for probiotic capsules versus yoghurts, possibly because

518    different strains and species are consumed with these products. Consumption of probiotic

519    capsules was associated with a lowering of the *Anaerococcus*-module, probiotics in general in

520    the last 24 hours with an increase of *L. gasseri* levels, while probiotic yoghurt decreased *L.*

521    *gasseri* levels. Unfortunately, our questionnaires lacked detailed information on the specific

522    species and strains in the probiotic products consumed by the Isala participants. Ultimately,

523    dedicated intervention studies with specific foods or diets, hygienic measures and/or

524    probiotic species and strains should further substantiate the associations found here, and

525    help the design of dedicated pharmaceutical and microbiome interventions.

526    **Conclusion**

527    In this large-scale remote-sampling study, we showed that the vaginal microbiome of women

528    from Belgium is mainly dominated by lactobacilli. We demonstrated that the vaginal

529    microbiome is a continuum, where taxon compositions that are in-between classical

530    community state types are frequently observed. Furthermore, we showed that most vaginal

531    taxa show small to moderate positive or negative abundance correlations with other taxa,

532    and that positively interacting vaginal taxa can be summarized by grouping them into modules

533    of intercorrelated taxa. In addition, we measured 166 participant covariates through

534    questionnaires. Our results showed that some of these factors explain a small but significant

535    part of vaginal microbiome variation, with "having had children" explaining the largest

536    fraction of the variation, after age. Finally, we highlighted that given conscious

537    communication tools and style, women are eager to participate in taboo-breaking

538    conversations as well as scientific studies aimed at improving their health. We therefore

539    endorse citizen science as a powerful approach to facilitate large-scale intimate microbiome

540    research and to empower citizens to impact their individual and community-level health by

541    promoting open science-based communication on taboo subjects.

542    **Acknowledgements**

558     ambassador Evi Hanssen, collaborative influencers on social media and supportive science

559     journalists for spreading the word on Isala and helping us build an online community that

560     openly discusses vaginal health with the aim to break the taboo.

561 **Author contributions**

562     SL, SA, EO, SW, GD, VV and CDB designed the study and worked on the conceptualization of

563     the research project. SL, SA, TG, TE, JD, SC, EO, IS, SW, CM and WVB worked on the

564     questionnaire set-up and cleaned the answers. SA, SL, JD, EO, TE and LVD carried out the

565     experimental and logistical work. SW and TG processed the sequencing data and performed

566     the biostatistical analyses. TG, SW, SA, SC and SL worked on the visualizations. SL, SW, TG, SA,

567     VV, GD, SC, JD, IS, PAB and CM contributed to the interpretation of the results. SL, SA, SW and

568     TG wrote the original manuscript. All authors contributed to reviewing and editing of the final

569     manuscript.

577 **Competing interests**

578     SL is a voluntary academic board member of ISAPP (the International Scientific Association on

579     Probiotics and Prebiotics, www.isappscience.org) and chairperson of the scientific advisory

580    board of YUN (yun.be). PAB is an independent consultant for several companies in the food

581    and pharmaceutical industry. GD is the chairperson of Femicare vzw (femicare.net) and has

582    worked as a medical consultant for various industries. However, none of these organizations

583    or companies was involved in the design, communication or data analysis of this Isala study,

584    which was fully funded by university, governmental and European funding, with the largest

585    part funded by the ERC StG project Lacto-Be.

586    **Methods**

587    <u>Study cohort and data collection</u>

588    The study was approved by the Ethical Committee of the Antwerp University

589    Hospital/University of Antwerp (B300201942076) and registered online at clinicaltrials.gov

590    with the unique identifier NCT04319536. The call for participants was launched on March

591    24[th], 2020 with the only inclusion criteria were being not pregnant and at least 18 years old.

592    Within ten days, 6,007 women registered through the Isala website (https://isala.be/en/) by

593    filling five questions on age, postal code, previous pregnancies, residence country in first

594    three years and contraceptive use. After obtaining a digital informed consent, these

595    participants were invited to fill out a large online questionnaire that included 137 relevant

596    and GDPR-compliant questions on the Qualtrics platform (Qualtrics, Provo, UT, USA). The

597    4,681 participants that filled out the entire questionnaire were invited to fill out their address

598    on the website to receive an Isala self-sampling kit. Eventually, 4,106 self-sampling kits were

599    sent out and 81.5% of the kits were returned to the University of Antwerp between July-

600    October 2020. Two vaginal swabs were self-collected in a standardized way by non-pregnant

601    participants (n = 3,323). And 3,294 participants filled out a short follow-up questionnaire with

602    39 questions within 24 hours of sampling.

603 Each kit contained two vaginal swabs. First the eNAT™ (Copan, Brescia, Italy), intended for

604 microbiome profiling, was collected and immediately afterwards the ESwab™ (Copan,

605 Brescia, Italy), intended for culturomics and metabolomics, was collected. In the insert it was

606 stipulated that both swabs had to be turned around 2-3 times to acquire enough biomass.

607 Immediately after sampling swabs were to be transferred to a vial which contained the

608 commercial transport buffer of the eNAT or ESwab and stored at home in the fridge. At last,

609 all samples were transported on room temperature with prepaid services by the national

610 parcel service (Bpost) with an average transport time of 2,9 +- 3,3 days (n = 3,306) from which

611 92,8% arrived within 7 days from sampling. Upon arrival, the eNAT swabs were stored at -

612 20°C until further processing in the lab[78]. The ESwab was vortexed for 15 seconds and

613 separated in two aliquots of 500µL, the first of which was stored at -80°C in a 96 tube Micronic

614 plate with 500µL 50% glycerol, the other being centrifuged for 3 min at 13,000 g, and its

615 supernatant stored in a 96 tube Micronic plate at -80°C as well.

616 <u>16S rRNA amplicon sequencing</u>

617 Before further processing, all samples were vortexed for 15-30 seconds and extracted with

618 the DNeasy PowerSoil Pro Kit of which some manually and other automated with the QIAcube

619 (Qiagen, Hilden, Germany) according to the instructions of the manufacturer. DNA

620 concentration of all samples was measured using the Qubit 3.0 Fluorometer (Life

621 Technologies, Ledeberg, Belgium) according to the instructions of the manufacturer. No less

622 than 2 µl of each bacterial DNA sample was used to amplify the V4 region of the 16S rRNA

623 gene, using standard barcoded forward (515F) and reverse (806R) primers[78]. These primers

624 were altered for dual index paired-end sequencing, as described in Kozich *et al.* (2013)[79]. The

625 resulting PCR products were checked on a 1.2% agarose gel. The PCR products were then

626    purified using the Agencourt AMPure XP Magnetic BeadCapture Kit (Beckman Coulter,

627    Suarlee, Belgium) and the concentration of all samples was measured using the Qubit 3.0

628    Fluorometer. Next, a library was prepared by pooling all PCR samples in equimolar

629    concentrations. This library was loaded onto a 0.8% agarose gel and purified using the

630    NucleoSpin Gel and PCR clean-up (Macherey-Nagel). The final concentration of the library was

631    measured with the Qubit 3.0 Fluorometer. Afterwards the library was denatured with 0.2N

632    NaOH (Illumina, San Diego California United States), diluted to 6 pM and spiked with 10-15%

633    PhiX control DNA (Illumina). Finally, dual-index paired-end sequencing was performed on a

634    MiSeq Desktop sequencer (Illumina). All DNA samples as well as negative controls of both PCR

635    (PCR grade water) and the DNA extraction runs were included on the sequencing runs. In

636    total, samples were sequenced across nine different MiSeq runs.

637    In order to validate our amplicon sequencing pipeline, including *Lactobacillus* subgenus

638    classification, we sequenced samples from the Isala pilot study in Ahannach, Delanghe, *et al.*

639    (2021)[78] with both amplicon and shotgun sequencing. These samples were processed in the

640    same way as the Isala samples, except that the DNA extraction was performed with the

641    HostZERO Microbial DNA Kit (Zymo Research, California, United States). These samples were

642    sequenced across two different MiSeq sequencing runs.

643    <u>Metagenomic shotgun sequencing (Isala pilot study samples)</u>

644    For the metagenomic shotgun sequencing of samples from the Isala pilot study, library

645    preparation was performed using the Nextera™ DNA Flex Library Prep or Nextera™ XT DNA

646    Library Preparation kit (Illumina), according to the instructions of the manufacturer. For the

647    Nextera™ DNA Flex Library Prep, 2 – 30 µL DNA sample was used to obtain input DNA with a

648    start amount between 1 and 100 ng. For the Nextera™ XT DNA Library Preparation kit, 1 ng

649  DNA samples in 5 μL was used as input DNA. For both protocols, when the 1 ng input DNA

650  could not be obtained for a certain DNA sample, the library preparation was continued with

651  the highest available amount of input DNA. Pooling of the libraries was done individually using

652  the Qubit 3.0 Fluorometer. During library preparation, library quality was checked using the

653  5200 Fragment Analyzer System with Agilent High Sensitivity NGS Fragment Kit (DNF-474).

654  22μL NGS Diluent Marker solution was mixed with 2μL library and ran on the Fragment

655  Analyzer, according the instruction of the manufacturer. The NGS DNA Ladder was used as

656  standard. Finally, the library was sequenced on a MiSeq desktop sequencer. In total, shotgun

657  samples were sequenced on two MiSeq runs.

658  <u>Creation of custom taxonomic reference databases</u>

659  In order to increase taxonomic resolution for the genus *Lactobacillus*, the genus was split into

660  nine subgenera. These subgenera were defined in three steps. First, a maximum-likelihood

661  species phylogeny of the genus was constructed using amino acid sequences of 100 single-

662  copy core genes from representative genomes, using the software IQ-TREE[80]. Second, the

663  subgenera were manually defined as the minimum number of clades in the species phylogeny

664  that would be needed to discriminate the four major vaginal *Lactobacillus* species. Finally, the

665  subgenera were checked for monophyly against the species phylogeny of release 05-RS95 of

666  the Genome Taxonomy Database (GTDB)[31].

667  To be able to classify amplicon sequences to the *Lactobacillus* subgenera, a custom 16S rRNA

668  reference database was created. This was done by downloading 16S rRNA sequences

669  extracted from sequenced genomes from the GTDB (release 05-RS95) as well as the GTDB

670  taxonomy hierarchy. This dataset was subsetted to sequences of the family *Lactobacillaceae*

671  only, and the genus *Lactobacillus* in the taxonomy hierarchy was replaced by the respective

672    subgenera of the species. Finally, these files were converted into a DADA2-compatible

673    reference database.

674    To be able to validate our amplicon data processing pipeline, including classification to

675    *Lactobacillus* subgenera, we also created a custom reference database for the classification

676    of metagenomic shotgun sequencing data. This database was created from three pieces of

677    data: (1) representative genomes for all bacterial species, downloaded from release 05-RS95

678    of the GTDB, (2) the GTDB taxonomy hierarchy updated with the *Lactobacillus* subgenera, and

679    (3) version GRCh38 of the human genome, downloaded from NCBI RefSeq[81]. These files were

680    used to create a database in Kraken2-compatible format.

681    <u>Processing and quality control of amplicon sequencing data</u>

682    Quality control and processing of amplicon reads was performed with the R package DADA2,

683    version 1.6.0[82]. First, reads with more than two expected errors were removed (no trimming

684    was performed). Next, paired reads were merged; in this process, read pairs with one or more

685    sequence conflicts were removed. Chimeras were then detected and removed with the

686    removeBimeraDenovo function. The merged and denoised reads (amplicon sequence

687    variants or ASVs) were taxonomically annotated from the phylum to the genus level with the

688    assignTaxonomy function using the EzBioCloud reference 16S rRNA database[83]. Next, three

689    different reclassifications were performed. First, ASVs classified to the family

690    *Leuconostocaceae* were reclassified to the family *Lactobacillaceae* to be in line with the recent

691    taxonomic update[13]. Second, the *Lactobacillaceae* ASVs were reclassified on the genus level

692    to the new genera defined by Zheng et al. And finally, ASVs of the updated genus *Lactobacillus*

693    (previously known as the *Lactobacillus delbrueckii* group) were reclassified to nine different

694    subgenera that we manually defined based on the phylogeny of the genus.

695    Taxon and sample quality control was performed as follows. Non-bacterial ASVs (e.g.,

696    mitochondria and chloroplasts) and ASVs with a length greater than 260 bases were removed.

697    Quality control of the samples was based on normalized read concentrations, which were

698    calculated as follows. First, the total read count per sample was divided by the volume of that

699    sample added to the sequencing library of its MiSeq run (there were nine runs in total). Next,

700    these read concentrations were normalized by dividing them by the median read

701    concentration of their respective run. Samples were then filtered using two criteria: (1) the

702    normalized read concentration should be higher than 0.05 and (2) the read count of a sample

703    should be greater than 2,000.

704    The Isala pilot study samples were processed in the same way as described above, with the

705    following exceptions: (1) ASV classification was performed with a 16S rRNA reference

706    database constructed from version 05-RS95 of the GTDB, followed by reclassification of the

707    *Lactobacillus* ASVs only to the custom *Lactobacillus* subgenera; (2) sample quality control was

708    based on a minimum read count of 1,000 reads.

709    <u>Processing and quality control of metagenomic sequencing data (pilot study samples)</u>

710    Metagenomic shotgun sequenced samples from the Isala pilot study were processed as

711    follows. First, paired reads were filtered with the DADA2 R package, version 1.20.0[82], requiring

712    a minimum length of 50 bases, a maximum of two uncalled bases per read and a maximum

713    of two expected errors per read. Next, read pairs were classified from the phylum to the

714    species level with Kraken2[84], using a custom reference database designed to validate our

715    amplicon sequencing pipeline (including *Lactobacillus* subgenus classification). Based on the

716    read classifications against this custom database, a read count table was constructed where

717    the columns represent taxa and the rows represent samples. Taxa were either species or

718  higher-level taxa for reads that were unclassified at one or more ranks. Non-bacterial taxa

719  were removed from the data, as were samples with fewer than 500 bacterial reads.

720  All processing of amplicon and shotgun datasets was performed in R version 4.1.1[85], using the

721  tidyverse set of packages, version 1.3.0[86], and the in-house package tidyamplicons, version

722  0.2.1.

723  <u>Culture analyses</u>

724  Based on the questionnaire answers a selection of self-reported "healthy" women was made.

725  This selection took place during the course of the study, so it does not include all "healthy"

726  women and included 592 women with: no known infection at the moment of sampling; no

727  use of vaginal probiotics; no current smokers; good general health; no use of

728  antibiotics/antimycotics in the past three months; no vaginal douching; no overall vaginal

729  conditions. The 592 samples were located in the detailed inventory and retrieved from the

730  Micronic plate at -80°C. The individual tubes were gathered to avoid melting of other samples

731  to preserve optimal viability of the microorganisms. To obtain single colonies, 10 μL of each

732  sample was inoculated on a small Petri dish (10mL) with three types of growth media (MRS,

733  MRS + vancomycin, or Colombia blood, all BD Difco™) and grown for 24-48h at 37°C and 5%

734  $CO_2$. After 24h the plates were checked for colonies and if present one colony of each plate

735  was selected at random, resulting in a maximum of three isolates per participant. A part of

736  this colony was inoculated in 10 mL MRS broth and grown overnight in 37°C and 5% $CO_2$. Of

737  the overnight grown culture, 800 μL was mixed with 800 μL 50% glycerol in labelled cryovials

738  (Greiner Bio-one Cryo.S™) and stored in -80°C. At the same time, another part of the colony

739  was also used for colony polymerase chain reaction (colony PCR) for taxonomic identification

740  with 16S Sanger sequencing, using universal primers 27F and 1492R.

Contraceptives, menstrual cycle and hormonal levels

742 Upon sampling, participants indicated when their menstrual cycle began, and also the average

743 length of their cycle. Depending upon the contraceptive, we used this data to determine the

744 day in which they are in, and predicted the levels of endo and exogenous levels of estrogen

745 and progestin. Peri and post-menopausal women were excluded from this analysis.

746 Statistical analyses

747 t-SNE-embeddings were performed on the relative abundances per sample, using the Bray-

748 Curtis distance metric[87] to calculate distances within the t-SNE[33]. Samples were classified into

749 a "primary type" based on the most dominant taxa, except if that taxon occurred less than

750 200 times as the most dominant taxon, in which case it was classified into a type "other". To

751 determine correlations between the abundances of taxa across our samples, we used the

752 fastspar implementation of SparCC with 100,000 permutations. We calculated correlations

753 only between taxa which were present at some non-zero abundance in at least 100 samples.

754 We used the same correlation threshold of 0.3 as in the original SparCC manuscript[34]. Clusters

755 were identified with hierarchical clustering with single linkage. Eigentaxa, a summary score

756 for a given set of taxa, (determined by the modules identified in the taxa-taxa correlation

757 networks) were calculated by first CLR-transforming the relative abundance data, and taking

758 the first principle component of the taxa in each cluster. Eigentaxa were multiplied by the sign

759 of the correlation coefficient between the eigentaxa and a representative taxon for each

760 cluster: *Gardnerella, Prevotella* and *Limosilactobacillus* for the BV, AV and *Lactobacillus*

761 modules, respectively.

762 Associations between microbial community composition and the questionnaire were

763 performed with an Adonis test, as implemented in the vegan package in R. For each effect of

764     interest, we tested three models. 1) ~ e_i, 2) ~ e_t + e_i, and 3) ~ e_t + age + e_i , where e_t

765     are technical effects, e_i is the effect of interest. Technical effects used were identical across

766     all experiments, and consisted of sequencing run, normalized read concentration and library

767     size, which were found to be strongly associated with the principal component s of the

768     relative abundance. In order to optimize computational performance, initially 1,000

769     permutations were performed for each effect of interest. A total of 10,000 permutations were

770     performed only for those effects which had p-values equal to 0.001.

771     Associations between Shannon diversity and variable collected via the questionnaire were

772     performed with a multiple linear regression, with three different models, as in the Adonis

773     test, 1) Diversity ~ e_i, 2) Diversity ~ e_t + e_i, and 3) Diversity ~ e_t + age + e_i.

774     Associations between the relative abundance of specific taxa and the questionnaire were

775     done with a multiple linear regression, with a model CLR(RA_I) ~ e_t + e_i, where RA_i refers

776     to the relative abundance of a taxa of interest, and CLR refers to the centered log ratio[88].

777     Associations between assigned community types and the questionnaire were performed with

778     a logistic regression, where, for each pair of community types T_A and T_B, we tested the

779     following three models: 1) I_T ~ e_i, 2) I_T ~ e_t + e_i, and 3) T_I ~ e_t + age + e_i, where I_T

780     is an indicator function whereby: I_T = 0 if sample is in T_A else 1 if sample is in T_B. Results

781     in figure 5 show the results for model 3, except for age, in which the results for model 2 are

782     shown.

783     For the Adonis model analysis of total explained variance, we included all significant factors

784     in a factorial Adonis test (Factors included are shown in figure 5). In order to perform this,

785     missing values in the questions were encoded as separate categories.

786    All data handling and visualization was performed in python and R version 4.1.0[85] using the

787    tidyverse set of packages and the in-house developed package tidyamplicons

788    (github.com/Swittouck/tidyamplicons).

789    Data availability

790    Sequencing data are available at the European Nucleotide Archive (ENA) under bioproject

791    PRJEB50407.

792    **References**

793    1.    Weinstein, L., Bogin, M., Howard, J. H. & Finkelstone, B. B. A survey of the vaginal
794          flora at various ages, with special reference to the Döderlein bacillus. *Am. J. Obstet.*
795          *Gynecol.* **32**, 211–218 (1936).

796    2.    Lash, A. F. & Kaplan, B. A Study of Döderlein's Vaginal Bacillus. *Oxford Univ. Press* **38**,
797          333–340 (2021).

798    3.    Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc. Natl. Acad. Sci.*
799          *U. S. A.* **108**, 4680–4687 (2011).

800    4.    Lopes dos Santos Santiago, G. *et al.* Longitudinal qPCR Study of the Dynamics of L.
801          crispatus, L. iners, A. vaginae, (Sialidase Positive) G. vaginalis, and P. bivia in the
802          Vagina. *PLoS One* **7**, e45281 (2012).

803    5.    El Aila, N. A. *et al.* Identification and genotyping of bacteria from paired vaginal and
804          rectal samples from pregnant women indicates similarity between vaginal and rectal
805          microflora. *BMC Infect. Dis.* **9**, 167 (2009).

806    6.    Oerlemans, E. F. M. *et al.* The Dwindling Microbiota of Aerobic Vaginitis, an
807          Inflammatory State Enriched in Pathobionts with Limited TLR Stimulation. *Diagnostics*
808          **10**, 879 (2020).

809    7.    Donders, G. G. G. *et al.* Definition of a type of abnormal vaginal flora that is distinct
810          from bacterial vaginosis: Aerobic vaginitis. *BJOG An Int. J. Obstet. Gynaecol.* **109**, 34–
811          43 (2002).

812    8.    Gosmann, C. *et al.* Lactobacillus-Deficient Cervicovaginal Bacterial Communities Are
813          Associated with Increased HIV Acquisition in Young South African Women. *Immunity*
814          **46**, 29–37 (2017).

815    9.    McClelland, R. S. *et al.* Evaluation of the association between the concentrations of
816          key vaginal bacteria and the increased risk of HIV acquisition in African women from
817          five cohorts: a nested case-control study. *Lancet Infect. Dis.* **18**, 554–564 (2018).

818    10.   Lewis, F. M. T., Bernstein, K. T. & Aral, S. O. Vaginal microbiome and its relationship to

819      behavior, sexual health, and sexually transmitted diseases. *Obstet. Gynecol.* **129**, 643–
820      654 (2017).

821 11.   Campisciano, G. *et al.* Subclinical alteration of the cervical–vaginal microbiome in
822      women with idiopathic infertility. *J. Cell. Physiol.* **232**, 1681–1688 (2017).

823 12.   Kroon, S. J., Ravel, J. & Huston, W. M. Cervicovaginal microbiota, women's health, and
824      reproductive outcomes. *Fertil. Steril.* **110**, 327–336 (2018).

825 13.   Zheng, J. *et al.* A taxonomic note on the genus Lactobacillus: Description of 23 novel
826      genera, emended description of the genus Lactobacillus Beijerinck 1901, and union of
827      Lactobacillaceae and Leuconostocaceae. *Int. J. Syst. Evol. Microbiol.* **70**, 2782–2858
828      (2020).

829 14.   Gajer, P. *et al.* Temporal Dynamics of the Human Vaginal Microbiota. *Sci. Transl. Med.*
830      **4**, 1–21 (2012).

831 15.   France, M. *et al.* VALENCIA: A Nearest Centroid Classification Method for Vaginal
832      Microbial Communities Based on Composition  1–15 (2020)
833      doi:10.21203/rs.2.24139/v1.

834 16.   Drell, T. *et al.* Characterization of the Vaginal Micro- and Mycobiome in
835      Asymptomatic Reproductive-Age Estonian Women. *PLoS One* **8**, (2013).

836 17.   Freitas, A. C. *et al.* The vaginal microbiome of pregnant women is less rich and
837      diverse, with lower prevalence of Mollicutes, compared to non-pregnant women. *Sci.*
838      *Rep.* **7**, 1–16 (2017).

839 18.   Lennard, K. *et al.* Microbial Composition Predicts Genital Tract Inflammation and
840      Persistent Bacterial Vaginosis in South African Adolescent Females. *Infect. Immun.* **86**,
841      (2017).

842 19.   Rhoades, N. S. *et al.* Longitudinal Profiling of the Macaque Vaginal Microbiome
843      Reveals Similarities to Diverse Human Vaginal Communities. *mSystems* **6**, (2021).

844 20.   Miller, E. A., Beasley, D. A. E., Dunn, R. R. & Archie, E. A. Lactobacilli dominance and
845      vaginal pH: Why is the human vaginal microbiome unique? *Front. Microbiol.* **7**, 1–13
846      (2016).

847 21.   Yildirim, S. *et al.* Primate vaginal microbiomes exhibit species specificity without
848      universal Lactobacillus dominance. *ISME J.* **8**, 2431–2444 (2014).

849 22.   Mirmonsef, P. *et al.* Free glycogen in vaginal fluids is associated with Lactobacillus
850      colonization and low vaginal pH. *PLoS One* **9**, 26–29 (2014).

851 23.   Song, S. D. *et al.* Daily Vaginal Microbiota Fluctuations Associated with Natural
852      Hormonal Cycle, Contraceptives, Diet, and Exercise. *mSphere* **5**, 1–14 (2020).

853 24.   Foxman, B., Muraglia, R., Dietz, J. P., Sobel, J. D. & Wagner, J. Prevalence of recurrent
854      vulvovaginal candidiasis in 5 European countries and the United States: Results from
855      an internet panel survey. *J. Low. Genit. Tract Dis.* **17**, 340–345 (2013).

856 25.   Medina, M. & Castillo-Pino, E. An introduction to the epidemiology and burden of
857      urinary tract infections. *Ther. Adv. Urol.* **11**, 3–7 (2019).

858  26.  Serrano, M. G. *et al.* Racioethnic diversity in the dynamics of the vaginal microbiome
859       during pregnancy. *Nat. Med.* **25**, 1001–1011 (2019).

860  27.  Noppe, J. *et al.* Vlaamse Migratie- en integratiemonitor 2018. *Brussel Agentschap*
861       *Binnenl. Best.* 311 (2018).

862  28.  Vlaanderen, S. Bevolking onder de armoededrempel - Statistiek Vlaanderen.
863       https://www.statistiekvlaanderen.be/nl/bevolking-onder-de-armoededrempel.

864  29.  Johnson, J. S. *et al.* Evaluation of 16S rRNA gene sequencing for species and strain-
865       level microbiome analysis. *Nat. Commun.* **10**, 1–11 (2019).

866  30.  Putonti, C., Shapiro, J. W., Ene, A., Tsibere, O. & Wolfe, A. J. Comparative Genomic
867       Study of Lactobacillus jensenii and the. *Am. Soc. Microbiol.* **5**, 1–5 (2020).

868  31.  Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity through
869       a phylogenetically consistent, rank normalized and complete genome-based
870       taxonomy. *Nucleic Acids Res.* **202**, 1–10 (2021).

871  32.  Rocha, J. *et al.* Lactobacillus mulieris sp. nov., a new species of lactobacillus
872       delbrueckii group. *Int. J. Syst. Evol. Microbiol.* **70**, 1522–1527 (2020).

873  33.  van der Maaten, L. & Hinton, G. Visualizing Data using t-SNE. *J. Mach. Learn. Res.* **9**,
874       2579–2605 (2008).

875  34.  Watts, S. C., Ritchie, S. C., Inouye, M. & Holt, K. E. FastSpar: Rapid and scalable
876       correlation estimation for compositional data. *Bioinformatics* **35**, 1064–1066 (2019).

877  35.  Rizzo, A., Losacco, A. & Carratelli, C. R. Lactobacillus crispatus modulates epithelial
878       cell defense against Candida albicans through Toll-like receptors 2 and 4, interleukin 8
879       and human β-defensins 2 and 3. *Immunol. Lett.* **156**, 102–109 (2013).

880  36.  Ojala, T. *et al.* Comparative genomics of Lactobacillus crispatus suggests novel
881       mechanisms for the competitive exclusion of Gardnerella vaginalis. *BMC Genomics*
882       **15**, 1–21 (2014).

883  37.  van der Veer, C. *et al.* Comparative genomics of human Lactobacillus crispatus
884       isolates reveals genes for glycosylation and glycogen degradation: Implications for in
885       vivo dominance of the vaginal microbiota. *Microbiome* **7**, 1–14 (2019).

886  38.  Falony, G. *et al.* Population-level analysis of gut microbiome variation. *Science* **352**,
887       560–564 (2016).

888  39.  Peltola, T. & Arpin, I. Science for everybody?: Bridging the socio-economic gap in
889       urban biodiversity monitoring. *Citiz. Sci. Innov. Open Sci. Soc. Policy* 367–380 (2018).

890  40.  Law, E. *et al. The Science of Citizen Science*. (2017). doi:10.1145/3022198.3022652.

891  41.  Petrova, M. I., Reid, G., Vaneechoutte, M. & Lebeer, S. Lactobacillus iners : Friend or
892       Foe? *Trends Microbiol.* **25**, 182–191 (2017).

893  42.  France, M. T. *et al.* Complete Genome Sequences of Six Lactobacillus iners Strains
894       Isolated from the Human Vagina. *Microbiol. Resour. Announc.* **9**, 17–19 (2020).

895  43.  Castro, J., Machado, D. & Cerca, N. Unveiling the role of Gardnerella vaginalis in

896     polymicrobial Bacterial Vaginosis biofilms: the impact of other vaginal pathogens
897     living as neighbors. *ISME J.* **13**, 1306–1317 (2019).

898  44. Harwich, M. D. *et al.* Drawing the line between commensal and pathogenic
899     Gardnerella vaginalis through genome analysis and virulence studies. *BMC Genomics*
900     **11**, (2010).

901  45. Łaniewski, P. & Herbst-Kralovetz, M. M. Bacterial vaginosis and health-associated
902     bacteria modulate the immunometabolic landscape in 3D model of human cervix. *npj*
903     *Biofilms Microbiomes* **7**, 1–17 (2021).

904  46. Castro, J., Jefferson, K. K. & Cerca, N. Genetic Heterogeneity and Taxonomic Diversity
905     among Gardnerella Species. *Trends Microbiol.* **28**, 202–211 (2020).

906  47. Charbonneau, M. R. *et al.* A microbial perspective of human developmental biology.
907     doi:10.1038/nature18845.

908  48. Koren, O. *et al.* A Guide to Enterotypes across the Human Body: Meta-Analysis of
909     Microbial Community Structures in Human Microbiome Datasets. *PLOS Comput. Biol.*
910     **9**, e1002863 (2013).

911  49. Canon, F., Nidelet, T., Guédon, E., Thierry, A. & Gagnaire, V. Understanding the
912     Mechanisms of Positive Microbial Interactions That Benefit Lactic Acid Bacteria Co-
913     cultures. *Front. Microbiol.* **11**, 1–16 (2020).

914  50. Blasche, S. *et al.* Metabolic cooperation and spatiotemporal niche partitioning in a
915     kefir microbial community. *Nat. Microbiol. |* **6**,.

916  51. Agarwal, K. *et al.* Glycan cross-feeding supports mutualism between Fusobacterium
917     and the vaginal microbiota. *PLOS Biol.* **18**, e3000788 (2020).

918  52. Mokoena, M. P. Lactic Acid Bacteria and Their Bacteriocins: Classification,
919     Biosynthesis and Applications against Uropathogens: A Mini-Review. *Mol.  A J. Synth.*
920     *Chem. Nat. Prod. Chem.* **22**, (2017).

921  53. Petrova, M. I., Lievens, E., Malik, S., Imholz, N. & Lebeer, S. Lactobacillus species as
922     biomarkers and agents that can promote various aspects of vaginal health. *Front.*
923     *Physiol.* **6**, (2015).

924  54. Lin, X. B. *et al.* The evolution of ecological facilitation within mixed-species biofilms in
925     the mouse gastrointestinal tract. *ISME J.* **12**, 2770–2784 (2018).

926  55. Hummelen, R. *et al.* Lactobacillus rhamnosus GR-1 and L. reuteri RC-14 to prevent or
927     cure bacterial vaginosis among women with HIV. *Int. J. Gynecol. Obstet.* **111**, 245–248
928     (2010).

929  56. Liu, J. J., Reid, G., Jiang, Y., Turner, M. S. & Tsai, C. C. Activity of HIV entry and fusion
930     inhibitors expressed by the human vaginal colonizing probiotic Lactobacillus reuteri
931     RC-14. *Cell. Microbiol.* **9**, 120–130 (2007).

932  57. Martinez, R. C. R. *et al.* Improved cure of bacterial vaginosis with single dose of
933     tinidazole (2 g), Lactobacillus rhamnosus GR-1, and Lactobacillus reuteri RC-14: A
934     randomized, double-blind, placebo-controlled trial. *Can. J. Microbiol.* **55**, 133–138

935          (2009).

936    58.   Verhelst, R. *et al. Cloning of 16S rRNA genes amplified from normal and disturbed*
937          *vaginal microflora suggests a strong association between Atopobium vaginae,*
938          *Gardnerella vaginalis and bacterial vaginosis*. http://www.biomedcentral.com/1471-
939          2180/4/16 (2004).

940    59.   Hardy, L. *et al.* A fruitful alliance: the synergy between Atopobium vaginae and
941          Gardnerella vaginalis in bacterial vaginosis-associated biofilm. *Sex. Transm. Infect.* **92**,
942          487–491 (2016).

943    60.   Randis, T. M. & Ratner, A. J. Gardnerella and Prevotella: Co-conspirators in the
944          Pathogenesis of Bacterial Vaginosis. *J. Infect. Dis.* **220**, 1085–1088 (2019).

945    61.   Donders, G. G. G., Bellen, G., Grinceviciene, S., Ruban, K. & Vieira-Baptista, P. Aerobic
946          vaginitis: no longer a stranger. *Res. Microbiol.* **168**, 845–858 (2017).

947    62.   Perrotta, A. R. *et al.* The Vaginal Microbiome as a Tool to Predict rASRM Stage of
948          Disease in Endometriosis: a Pilot Study. doi:10.1007/s43032-019-00113-5.

949    63.   Haggerty, C. L. *et al.* Presence and concentrations of select bacterial vaginosis-
950          associated bacteria are associated with increased risk of pelvic inflammatory disease.
951          *Sex. Transm. Dis.* **47**, 344 (2020).

952    64.   De Seta, F., Campisciano, G., Zanotta, N., Ricci, G. & Comar, M. The vaginal
953          community state types microbiome-immune network as key factor for bacterial
954          vaginosis and aerobic vaginitis. *Front. Microbiol.* **10**, 2451 (2019).

955    65.   Si, J., You, H. J., Yu, J., Sung, J. & Ko, G. P. Prevotella as a Hub for Vaginal Microbiota
956          under the Influence of Host Genetics and Their Association with Obesity. *Cell Host*
957          *Microbe* **21**, 97–105 (2017).

958    66.   Vodstrcil, L. A. *et al.* Combined oral contraceptive pill-exposure alone does not reduce
959          the risk of bacterial vaginosis recurrence in a pilot randomised controlled trial. *Sci.*
960          *Rep.* **9**, 1–13 (2019).

961    67.   Nelson, T. M. *et al.* Cigarette smoking is associated with an altered vaginal tract
962          metabolomic profile. *Sci. Rep.* **8**, 852 (2018).

963    68.   Lewis, C. A. *et al.* Effects of Hormonal Contraceptives on Mood: A Focus on Emotion
964          Recognition and Reactivity, Reward Processing, and Stress Response. *Curr. Psychiatry*
965          *Rep.* **21**, 1–15 (2019).

966    69.   Lundin, C., Wikman, A., Bixo, M., Gemzell-Danielsson, K. & Sundström Poromaa, I.
967          Towards individualised contraceptive counselling: clinical and reproductive factors
968          associated with self-reported hormonal contraceptive-induced adverse mood
969          symptoms. *BMJ Sex. Reprod. Heal.* **47**, e1–e8 (2021).

970    70.   Burrows, L. J., Basha, M. & Goldstein, A. T. The Effects of Hormonal Contraceptives on
971          Female Sexuality: A Review. *J. Sex. Med.* **9**, 2213–2223 (2012).

972    71.   Khialani, D., Rosendaal, F. & Vlieg, A. V. H. Hormonal Contraceptives and the Risk of
973          Venous Thrombosis. *Semin. Thromb. Hemost.* **46**, 865–871 (2020).
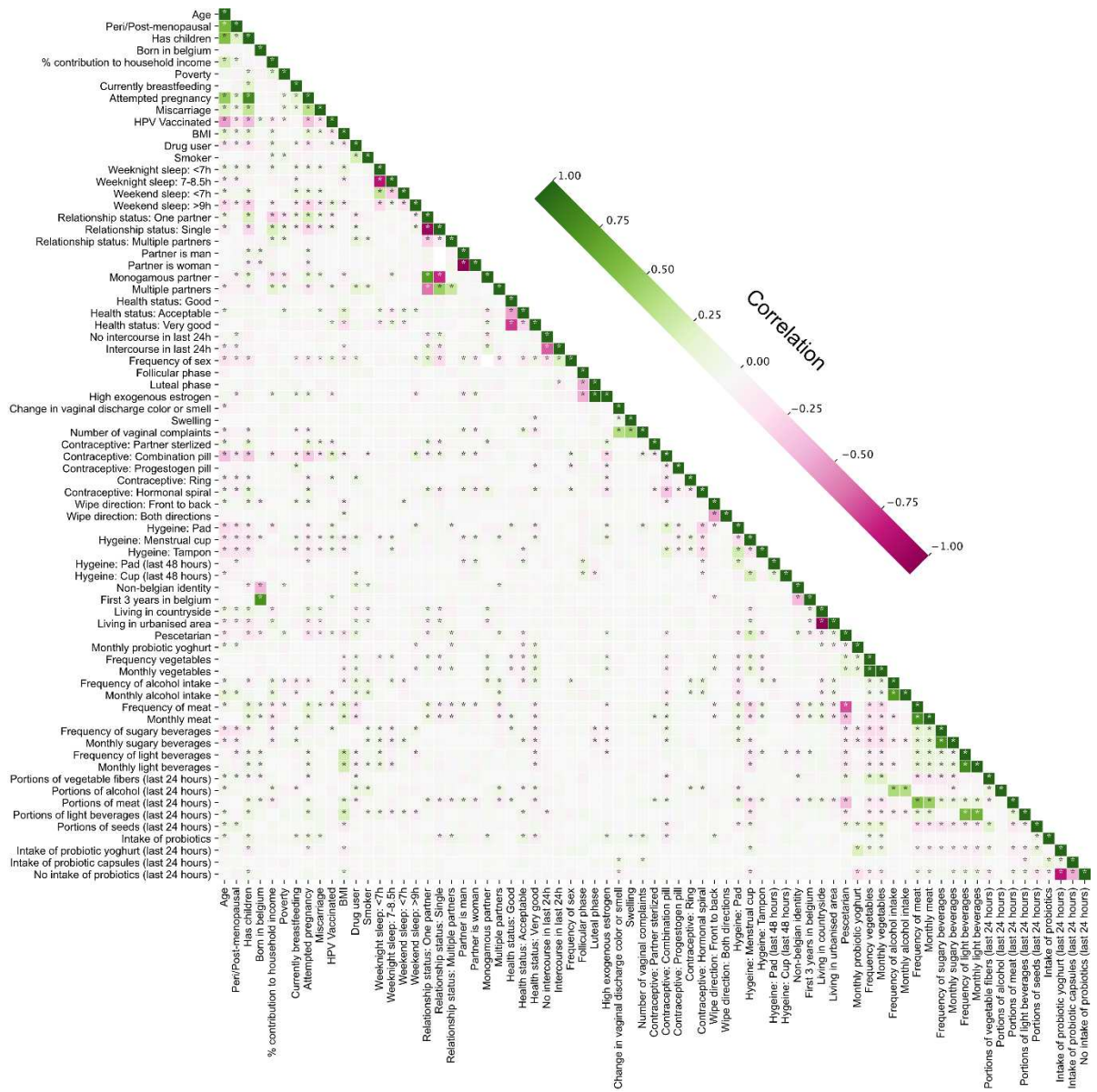
974    72.    Morimont, L., Haguet, H., Dogné, J. M., Gaspard, U. & Douxfils, J. Combined Oral
975          Contraceptives and Venous Thromboembolism: Review and Perspective to Mitigate
976          the Risk. *Front. Endocrinol. (Lausanne).* **12**, 1 (2021).

977    73.    Donders, G. G. G. *et al.* Screening for abnormal vaginal microflora by self-assessed
978          vaginal pH does not enable detection of sexually transmitted infections in Ugandan
979          women. *Diagn. Microbiol. Infect. Dis.* **85**, 227–230 (2016).

980    74.    Achilles, S. L., Meyn, L. A., Austin, M. N., Avolia, H. A. & Hillier, S. L. A longitudinal
981          evaluation of the impact of contraceptive initiation on vaginal microbiota in us
982          women. *Am. J. Obstet. Gynecol.* **219**, 643–644 (2018).

983    75.    Ferretti, P. *et al.* Mother-to-Infant Microbial Transmission from Different Body Sites
984          Shapes the Developing Infant Gut Microbiome. *Cell Host Microbe* **24**, 133-145.e5
985          (2018).

986    76.    DiGiulio, D. B. *et al.* Temporal and spatial variation of the human microbiota during
987          pregnancy. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 11060–11065 (2015).

988    77.    Lee, E. & Lee, J. E. Impact of drinking alcohol on gut microbiota: recent perspectives
989          on ethanol and alcoholic beverage. *Curr. Opin. Food Sci.* **37**, 91–97 (2021).

990    78.    Ahannach, S. *et al.* Microbial enrichment and storage for metagenomics of vaginal,
991          skin, and saliva samples. *iScience* **24**, 103306 (2021).

992    79.    Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K. & Schloss, P. D.
993          Development of a dual-index sequencing strategy and curation pipeline for analyzing
994          amplicon sequence data on the MiSeq Illumina sequencing platform. *Appl. Environ.*
995          *Microbiol.* **79**, 5112–20 (2013).

996    80.    Nguyen, L. T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A fast and
997          effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol.*
998          *Biol. Evol.* **32**, 268–274 (2015).

999    81.    O'Leary, N. A. *et al.* Reference sequence (RefSeq) database at NCBI: Current status,
1000         taxonomic expansion, and functional annotation. *Nucleic Acids Res.* **44**, D733–D745
1001         (2016).

1002    82.    Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina amplicon
1003         data. *Nat. Methods* **13**, 581–583 (2016).

1004    83.    Yoon, S. H. *et al.* Introducing EzBioCloud: A taxonomically united database of 16S
1005         rRNA gene sequences and whole-genome assemblies. *Int. J. Syst. Evol. Microbiol.* **67**,
1006         1613–1617 (2017).

1007    84.    Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification
1008         using exact alignments. *Genome Biol.* **15**, R46 (2014).

1009    85.    R Core Team. R: A Language and Environment for Statistical Computing. (2020).

1010    86.    Wickham, H. *et al.* Welcome to the Tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).

1011    87.    van der Maaten, L. Barnes-Hut-SNE. *1st Int. Conf. Learn. Represent. ICLR 2013 - Conf.*
1012         *Track Proc.* 1–11 (2013).

1013    88.    Aitchison, J. A concise guide to compositional data analysis. in *2nd Compositional*
1014           *Data Analysis Workshop* (2003).
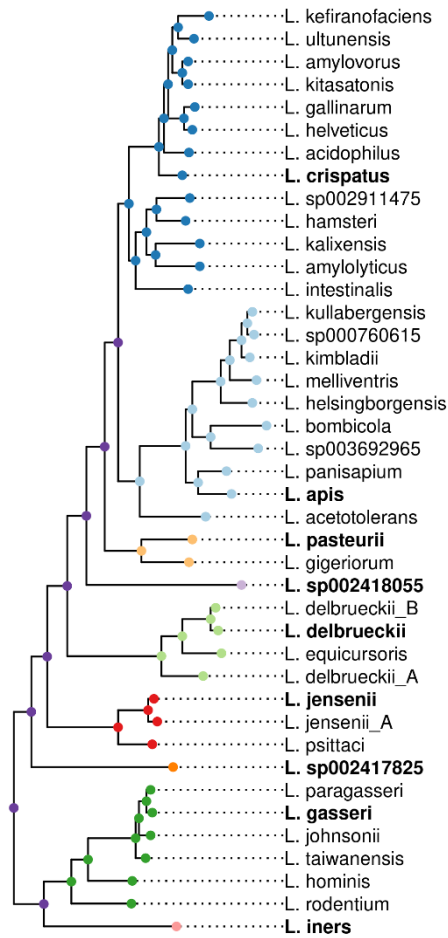
1015

1016
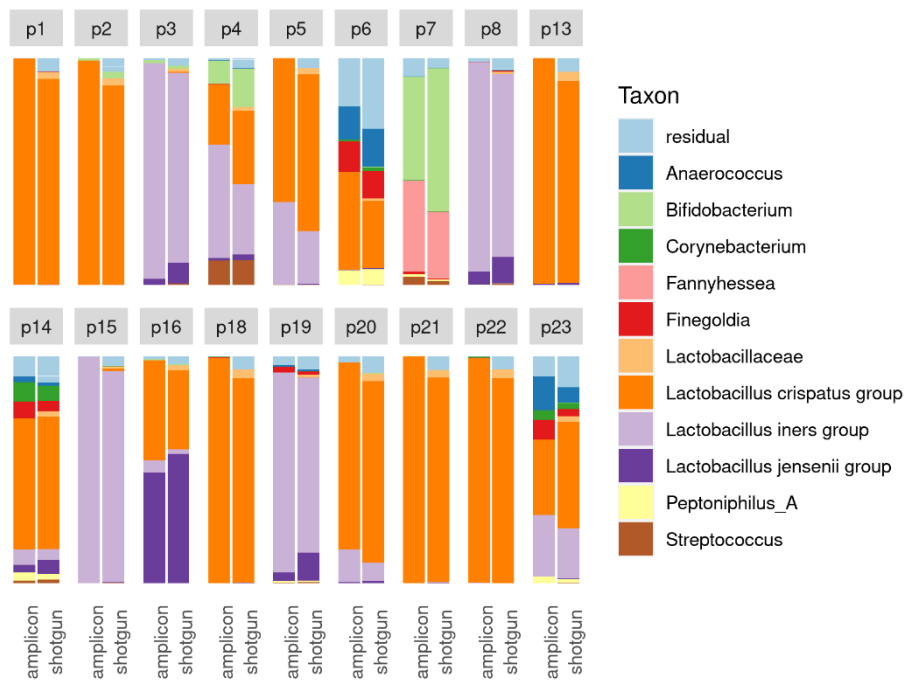
1017    **Supplementary figures**



1018

1019    Supplementary Figure 1 - **Correlations between a subset of the questionnaire variables.** Heatmap
1020    with correlations between questionnaire variables shown in Figure 5. Positive correlations are
1021    indicated with green, negative correlations in red. Significant correlations are marked with an asterisk.

1022

1023 Supplementary Figure 2 - **Species tree of *Lactobacillus* from the Genome Taxonomy Database.**
1024 Maximum-likelihood species phylogeny of the genus *Lactobacillus*, obtained by taking a subtree of the
1025 species phylogeny of the domain Bacteria inferred by the Genome Taxonomy Database (GTDB),
1026 release 05-RS95[31]. Colors indicate the nine custom-defined subgenera used in this study. Bold tip
1027 labels indicate representative species of the subgenera.

1028

1029 Supplementary Figure 3 - **Comparison between amplicon and shotgun sequencing results for 18**
1030 **samples.** Relative abundances for the eleven most abundant taxa overall. Each facet shows a vaginal
1031 sample from a single participant, sequenced with 16S rRNA amplicon sequencing (left) or
1032 metagenome shotgun sequencing (right).

1033

1034    Supplementary Figure 4 - **Example of a personal vaginal microbiome profile result.** Top left figure
1035    indicates the dominant type. Bottom left show the percentage ("verdeling") of the top eight taxa
1036    identified. Right figure (pie chart) displays the top six taxa plus the remaining ("overage") ones.



1037

**1398**

With Isala, we have found that this bacterium was dominant in the vagina of 1,398 women. That is about 43% of all participants that donated a sample.

## What does this bacterium look like?

Lactobacillus crispatus is a fairly long rod of 2 to 11 micrometers in size with a thick wall. That's not that big when you know that 1000 micrometers fit into one millimeter. The name comes from the English 'curled, crisped'. This bacterium was first discovered by Brygoo and Aladame in 1953.

## What does science already know about this bacterium?

Kind of a lot! This bacterium has a very extensive genome of about 2 million base pairs with more than 2000 genes, which means that this bacterium can make more than 2000 different proteins. She also seems to be well equipped to survive in a relatively wide variety of animal and human environments.

## What is this bacterium doing in my vagina?

Lactobacillus crispatus is very often associated with a healthy vagina. This bacterium produces a lot of lactic acid and therefore ensures acidity in the vagina. In this way, this bacterium protects your vagina against infections or pathogenic bacteria and fungi. Lactobacillus crispatus also makes other molecules that act as natural antibiotics or protect against inflammation, but not all these molecules are well known. When researching a healthy vaginal microbiome, we often focus on lactic acid, but each strain of Lactobacillus also produces an array of protective or beneficial molecules for our health.

Unravelling these molecules is something that Isala's team is happy to work on in the future. For example, we already know that Lactobacillus crispatus has a very good and active immune system so that this bacterium can protect itself against bacteriophages. These are viruses that can make (healthy) bacteria sick.

## Does this bacterium occur elsewhere?

Yes, Lactobacillus crispatus is also found in your gut and scientists have also found it in chickens. If you enter this bacterium in a search engine on the internet, you will probably come across a number of probiotics. After all, a lot of scientific research has already been done into the health effects of this bacterium.
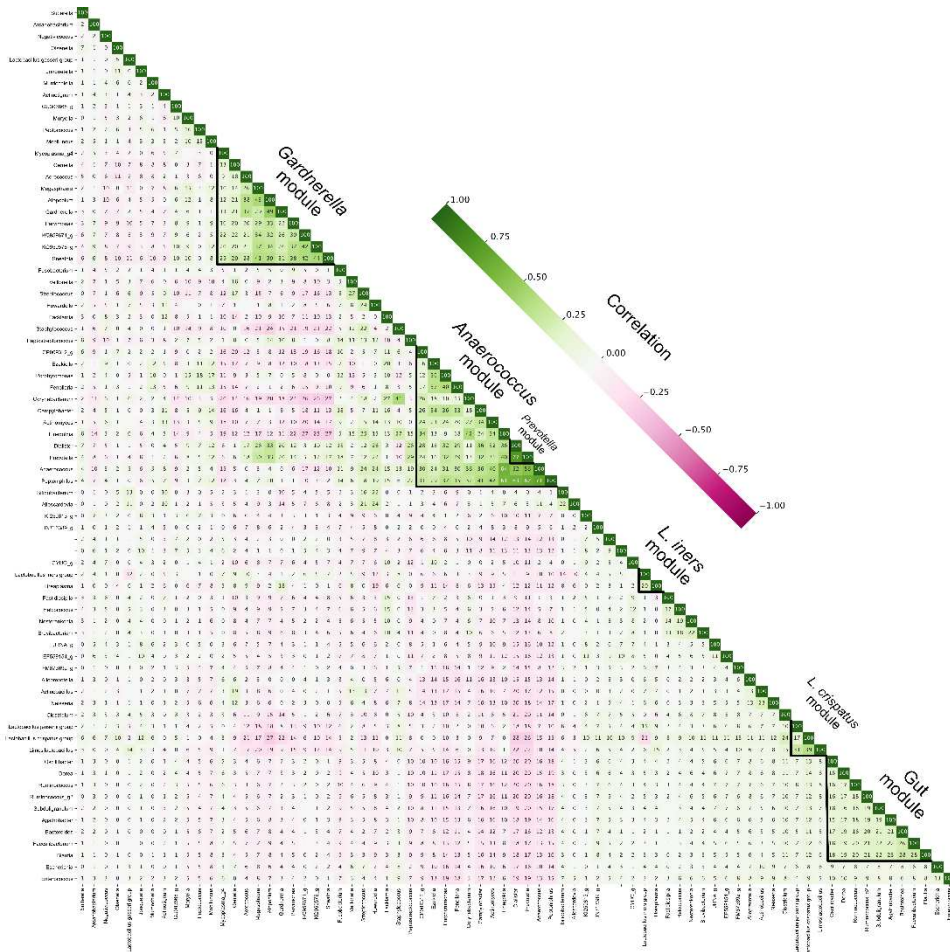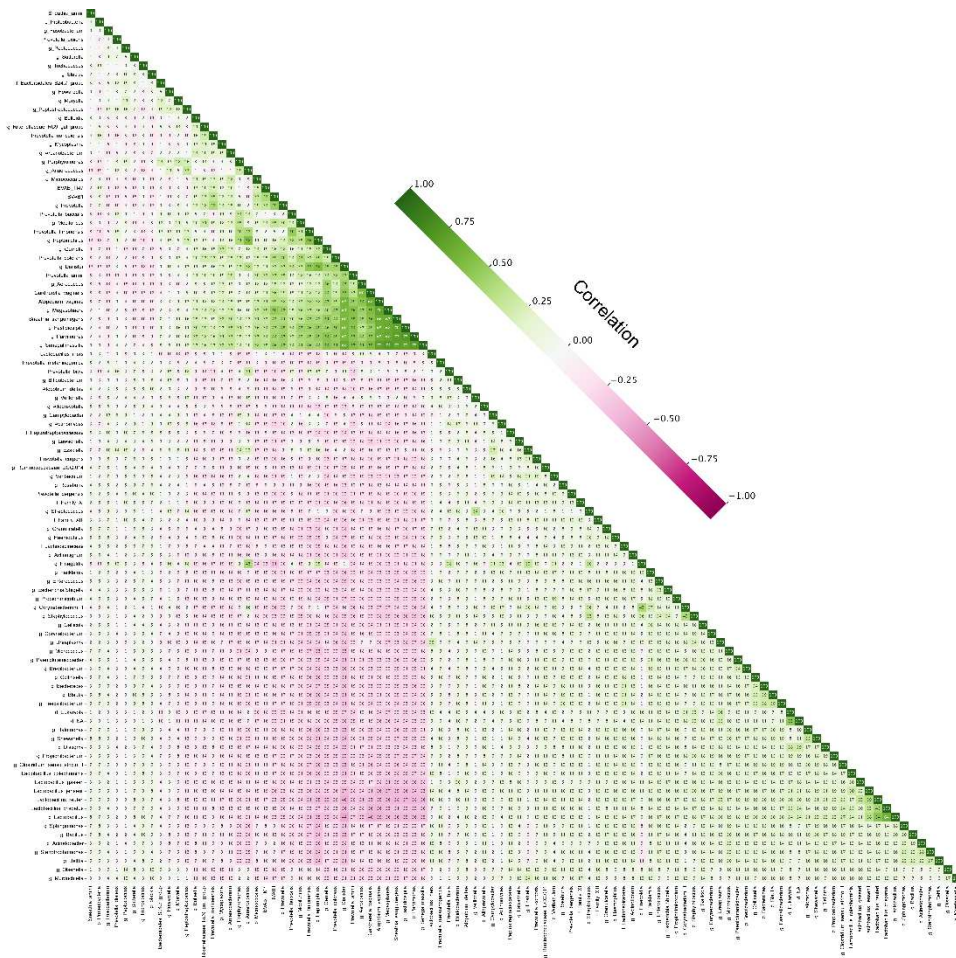
1038

1039 Supplementary Figure 5 - **Received information for non-microbiology experts.** To each of the top
1040 eight taxa a webpage was dedicated. Here, an example of the page on *Lactobacillus crispatus* is added.
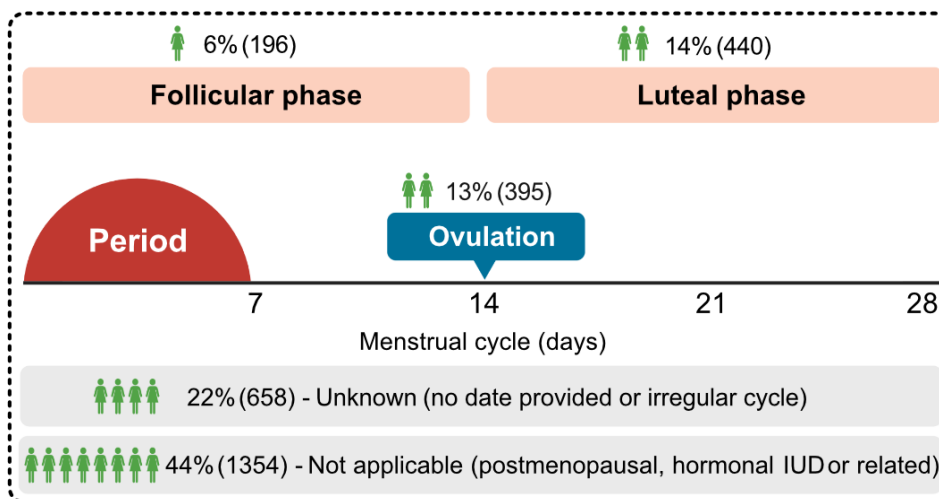1041 Other taxa can be accessed via https://isala.be/en/category/vaginal-bacteria/

1042

1043

Supplementary Figure 6 – **Full taxon correlation matrix.** SparCC correlation Network between taxa
determined in the Isala data. Positive correlations are indicated in green, negative correlations in red.
In each cell is given the correlation (*100) for each pair of taxa. The modules identified and shown in
figure 4 are indicated with triangles in the figure.

1048

Supplementary Figure 7 – **Full taxon correlation matrix of the Valencia study.** SparCC correlation Network between taxa determined in the Isala data. Positive correlations are indicated in green, negative correlations in red. In each cell is given the correlation (*100) for each pair of taxa.



1052

Supplementary Figure 8 – **The menstrual cycle.** Using information about each participant's cycle length and last menstruation, we estimated the stage of the cycle in which the swab was sampled. Participants whose cycles had irregular lengths, or who did not report their last menstruation were classified as "Unknown", and participants using hormonal contraceptives or were peri/post-menopausal were classified as "Not applicable".

1059 Supplementary Table 1 – **ASV occurrence and abundance of top 10 lactobacilli and**
1060 **percentage of top 10 isolated lactobacilli from Isala's samples.** The occurrence of the top 10
1061 ASVs of lactobacilli on (sub)genus level over all Isala's samples and their mean relative abundance and
1062 the percentage of isolates belonging to the top 10 most isolated lactobacilli (determined by 16S
1063 amplicon sequencing) in relation to the total lactobacilli isolates (n = 230) and the total number of
1064 isolates per species.

| (sub)genus | | | 16S isolates | | |
|---|---|---|---|---|---|
| (Sub)genus | Occurrence | Mean relative abundance | Species | Percentage of total lactobacilli isolates on De Man, Rogosa en Sharpe or Columbia Blood media (glucose as main sugar) | Number of isolates studied (n = 230) |
| *Lactobacillus crispatus group* | 0,897699005 | 0,399114797 | *Limosilactobacillus fermentum* | 24,49% | 60 |
| *Lactobacillus iners group* | 0,719527363 | 0,240823923 | *Lactobacillus crispatus* | 13,88% | 34 |
| *Limosilactobacillus* | 0,478544776 | 0,004111909 | *Lactobacillus jensenii* | 12,24% | 30 |
| *Lactobacillus jensenii group* | 0,467661692 | 0,04856063 | *Lactobacillus paragasseri* | 9,80% | 24 |
| *Lactobacillus gasseri group* | 0,268345771 | 0,029760051 | *Lacticaseibacillus rhamnosus* | 8,98% | 22 |
| *Lactobacillaceae* | 0,027052239 | 0,00199087 | *Lacticaseibacillus paracasei* | 7,35% | 18 |
| *Lacticaseibacillus* | 0,023942786 | 0,000494929 | *Limosilactobacillus reuteri* | 6,12% | 15 |
| *Lactiplantibacillus* | 0,00528607 | 1,65417E-05 | *Lactiplantibacillus plantarum* | 4,90% | 12 |
| *Ligilactobacillus* | 0,004975124 | 2,75966E-05 | *Lactobacillus gasseri* | 3,67% | 9 |
| *Apilactobacillus* | 0,003109453 | 0,000293372 | *Leuconostoc mesenteroides* | 2,45% | 6 |

1065

1066 Supplementary Table 2 – **Descriptive statistics of taxa.** Various descriptive statistics for
1067 subgenera of the genus *Lactobacillus* and genera detected in this study: number of ASVs
1068 within the (sub)genus (n_asvs), occurrence, average relative abundance
1069 (mean_rel_abundance), frequency of being the most abundant taxon and greater than 0%
1070 abundant (top_and_gt0p), same as previous but greater than 30% abundant
1071 (top_and_gt30p), same as previous but greater than 50% abundant (top_and_gt50p), the
1072 previous three measures but in terms of relative frequencies (top_and_gtXp_rel).

1073 Supplementary Table 3 – **Association tests between participant characteristics and their**
1074 **vaginal microbiome.** Results of statistical tests for each tested questionnaire responses.
1075 Results are provided for the beta-diversity (Adonis), alpha-diversity, taxa relative abundances
1076 and eigentaxa level tests. In addition to effect sizes, confidence intervals and p-values the
1077 number of participants in each condition are provided.

# Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- TableS2genusstats.csv
- TableS3associationtests20211220.xlsx