

## Microdiversity of the Vaginal Microbiome is Associated with Preterm Birth

Jingqiu Liao<sup>1,2\*†</sup>, Liat Shenhav<sup>3\*</sup>, Myrna Serrano<sup>4,5</sup>, Bin Zhu<sup>4,5</sup>, Gregory A. Buck<sup>4,5,6</sup>, Tal Korem<sup>1,7,8†</sup>

<sup>1</sup> Program for Mathematical Genomics, Department of Systems Biology, Columbia University Irving Medical Center, New York, NY, USA

<sup>2</sup> Department of Civil and Environmental Engineering, Virginia Tech, Blacksburg, VA, USA

<sup>3</sup> Center for Studies in Physics and Biology, Rockefeller University, New York, NY, USA

<sup>4</sup> Department of Microbiology and Immunology, School of Medicine, Virginia Commonwealth University, Richmond, VA, USA

<sup>5</sup> Center for Microbiome Engineering and Data Analysis, Virginia Commonwealth University, Richmond, VA, USA

<sup>6</sup> Department of Computer Science, School of Engineering, Virginia Commonwealth University, Richmond, VA, USA

<sup>7</sup> Department of Obstetrics and Gynecology, Columbia University Irving Medical Center, New York, NY, USA

<sup>8</sup> CIFAR Azrieli Global Scholars program, CIFAR, Toronto, Canada

\* These authors contributed equally to this work.

† Corresponding authors: Emails: [liaoqj@vt.edu](mailto:liaoqj@vt.edu), [tal.korem@columbia.edu](mailto:tal.korem@columbia.edu)

### Abstract

Preterm birth (PTB) is the leading cause of neonatal morbidity and mortality. The vaginal microbiome has been associated with PTB, yet the mechanisms underlying this association are not fully understood. Understanding microbial genetic adaptations to selective pressures, especially those related to the host, may yield new insights into these associations. To this end, we analyzed metagenomic data from 705 vaginal samples collected longitudinally during pregnancy from 40 women who delivered preterm spontaneously and 135 term controls from the Multi-Omic Microbiome Study-Pregnancy Initiative (MOMS-PI<sup>1</sup>). We find that the vaginal microbiome of pregnancies that ended preterm exhibits unique genetic profiles. It is more genetically diverse at the species level and harbors a higher richness and diversity of antimicrobial resistance genes, likely promoted by transduction. Interestingly, we find that *Gardnerella* species, a group of central vaginal pathobionts, are driving this higher genetic diversity, particularly during the first half of the pregnancy. We further present evidence that *Gardnerella* spp. undergoes more frequent recombination and stronger purifying selection in genes involved in lipid metabolism. Overall, our results reveal novel associations between the vaginal microbiome and PTB using population genetics analyses, and suggest that evolutionary processes acting on the vaginal microbiome may play a vital role in adverse pregnancy outcomes such as preterm birth.

## Introduction

Preterm birth (PTB), childbirth at <37 weeks of gestation, is the leading cause of neonatal morbidity and mortality<sup>2</sup>. Each year, approximately 15 million infants are born preterm globally, over 500,000 of them in the US<sup>3</sup>. Preterm infants are at a high risk of respiratory, gastrointestinal and neurodevelopmental complications<sup>4</sup>. While a number of maternal, fetal, and environmental factors have been associated with PTB<sup>2,5,6</sup>, its etiopathology remains largely unknown, and early diagnosis and effective therapeutics are still lacking.

Over the past decades, growing evidence has pointed to potential involvement of the vaginal microbiome in PTB<sup>1,7-10</sup>. This involvement has so far been mostly characterized as an ecological process, meaning changes in microbial abundances and vaginal community states. For instance, increased richness and diversity of microbial communities and the presence of particular community state types (CST), have been repeatedly associated with PTB<sup>1,9,11-15</sup>. In addition, vaginal microbiomes of women who delivered preterm appear to be less stable during pregnancy, with some studies reporting a significant decrease in the richness and diversity of these microbial communities during pregnancy<sup>1,12</sup>.

Multiple endogenous factors, such as hormonal changes, nutrient availability and microbial interactions, and exogenous factors, such as genital infections, antibiotic treatment and exposure to xenobiotics, could trigger ecological processes and alter the vaginal microbial composition<sup>16,17</sup>. These factors may also act as selective pressures that affect genetic variation in the microbial populations that make the vaginal microbiome. Such adaptive evolution in the vaginal environment, even during pregnancy, is highly plausible given the high mutation rates, short generation times, and large population sizes of microbes<sup>18</sup>. They are further supported by observations of rapid adaptation to environmental changes in other human-associated microbial ecosystems<sup>19-21</sup>. The way by which vaginal microbes respond to various selective pressures may, in turn, affect the host, including pregnancy outcomes. Therefore, a comprehensive investigation of the genetic diversity of the vaginal microbiome at the population level, which we term “microdiversity”, and the underlying evolutionary forces that shape it, holds promise for a better understanding of the etiopathology of PTB.

Here, we performed an in-depth population genetics analysis and characterized the population structure of the vaginal microbiome along pregnancy and in the context of preterm birth. We used metagenomic sequencing data from 705 vaginal samples collected longitudinally during pregnancy as part of the Multi-Omic Microbiome Study-Pregnancy Initiative (MOMS-PI<sup>1</sup>). Our analyses include samples from 40 women who subsequently experienced spontaneous preterm birth (sPTB) and 135 women who had a term birth (TB). We show that the vaginal microbiome of pregnancies that ended preterm exhibits higher nucleotide diversity at the species level and higher antimicrobial resistance potential. We find that this higher nucleotide diversity is driven by *Gardnerella* spp., a group of central vaginal pathobionts, especially during the first half of pregnancy, and suggest that this may be related to optimization of growth rates in this taxon. We further identify a strong association between evolutionary signatures and sPTB in *Gardnerella* spp., including more frequent homologous recombination and stronger purifying

selection. Overall, our results show novel associations between the vaginal microbiome and sPTB at the population genetics level, and suggest that evolutionary processes acting on the vaginal microbiome may play a critical role in sPTB, and potentially also in other adverse pregnancy outcomes.

## Results

### The phylogenetic composition of the vaginal microbiome associates with sPTB.

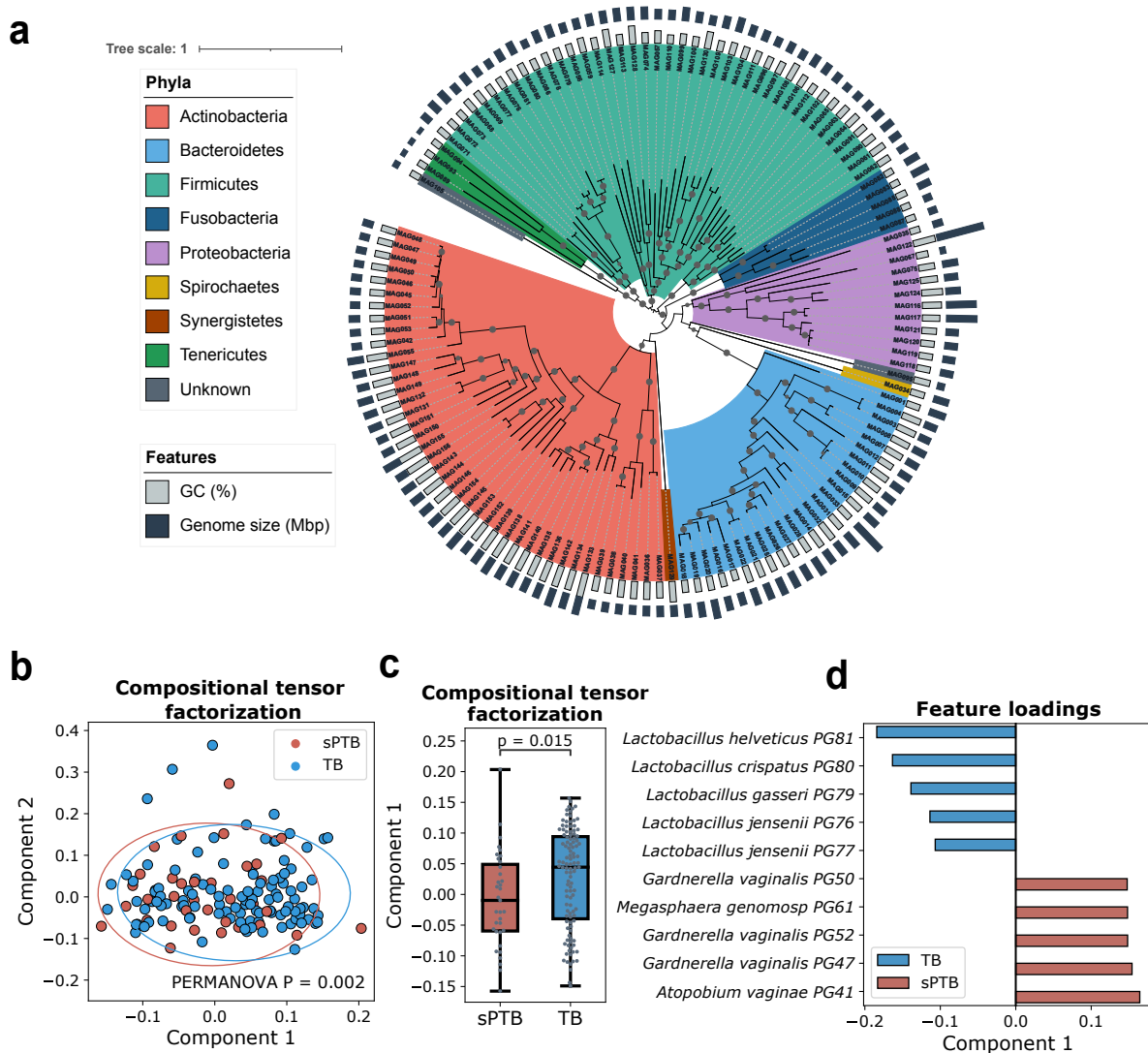
We assembled a total of 1,078 metagenome-assembled genomes (MAGs) with at least medium quality<sup>22</sup> (>50% completeness, <10% contamination; **Supplementary Table 1**; Methods) from previously published<sup>1</sup> metagenomic reads generated from 705 vaginal samples<sup>1</sup>. These samples were collected from 175 women visiting maternity clinics in Virginia and Washington at various time points along pregnancy<sup>1</sup>. We clustered these MAGs into 157 species-level phylogroups at the level of 95% average nucleotide identity (ANI), which roughly corresponds to the species level<sup>23</sup>; and selected the most complete MAG with the least contamination as the representative for each phylogroup. These representative MAGs were 86±14% (mean±SD) complete and 1.1±1.8% contaminated, with 93 (59% of 157) of them estimated to have high quality<sup>22</sup> (>90% completeness and <5% contamination; **Supplementary Table 1**). Taxonomic assignment of these representative MAGs (Methods) revealed that the phylogroups represent at least 8 phyla, with genome size (adjusted by completeness) ranging from 0.6 to 7.4 Mbps and GC content ranging from 25.3% to 69.7% (**Fig. 1a**, **Supplementary Table 1**). *Actinobacteria* had the most phylogroups detected in the samples, followed by *Firmicutes* and *Bacteroidetes* (**Fig. 1a**).

Of note, 12 of these species-level phylogroups (PG042-PG053) were assigned to *Gardnerella vaginalis* according to CheckM<sup>24</sup>, supporting the existence of multiple genotypes at the species level within the 'species' *G. vaginalis*<sup>25,26</sup>. To better resolve the classification of these *G. vaginalis* phylogroups, we compared the average nucleotide identity (ANI) for the representative MAGs of these phylogroups against updated reference genomes for *Gardnerella*, including *G. vaginalis*, *G. piotti*, *G. leopoldii*, *G. swidsinskii*, and nine species remained to be characterized (gs2-3 and gs7-13)<sup>25</sup>; gs-2-3 and gs7-13 correspond to group 2-3 and 7-13 shown in Fig. 1 in <sup>25</sup>. The ANI analysis shows that PG043 represents *G. vaginalis*, PG044 represents *G. swidsinskii*, PG042 represents *G. piotti*, and PG046, PG049, PG051 and PG053 represent *G. gs7*, *G. gs8*, *G. gs13* and *G. gs12*, respectively (**Supplementary Fig. 1**). The remaining phylogroups (PG045, PG047, PG048, PG050, and PG052) do not cluster with any reference species and may represent novel species of *Gardnerella*. Here, we refer to phylogroups PG042-PG053 as *Gardnerella* spp.

To understand if the temporal dynamics of the vaginal microbiome is associated with sPTB, we employed a revised version of compositional tensor factorization (CTF)<sup>27</sup> to assess temporal changes to the composition of the microbiome during pregnancy. This analysis shows a significant separation of women by pregnancy outcomes (PERMANOVA  $F = 8.0492$ ;  $P = 0.002$ ; **Fig. 1b**) based on the dynamics of their microbiome composition over time, specifically

observed for component 1 in the CTF analysis (Mann-Whitney  $P = 0.015$ ; **Fig. 1c**). We further found that the top features contributing to this difference belong to *Lactobacillus helveticus* (PG081), *Lactobacillus crispatus* (PG080), *Lactobacillus gasseri* (PG079), and *Lactobacillus jensenii* (PG076 and PG077) that are associated with TB and *Megasphaera genomsp.* (PG061), *Gardnerella* spp. (PG047, PG050, PG052), and *Atopobium vaginae* (PG041) that are associated with PTB (**Fig. 1d**); these species were previously found to be associated with pregnancy outcomes<sup>1,9,28</sup>. These results suggest that the vaginal microbiome has a different temporal trajectory during pregnancies ending preterm, consistent with previous findings<sup>1,7</sup>, and with *Gardnerella* as an important factor. Overall, our results demonstrate that de-novo metagenomic analysis replicates and expands previous findings with respect to associations between the composition of the vaginal microbiome and sPTB.

Next, we sought to examine the diversity of microbial strains detected within species, and its association with sPTB. We found strains of *M. genomsp.* showed significantly higher ANI between women who delivered preterm compared to a null distribution calculated based on ANIs from any two randomly selected women (Permutation  $P = 0.002$ , adjusted  $P < 0.05$ ; Methods), a relationship not observed between women who delivered at term ( $P = 0.208$ ; **Supplementary Fig. 2**). This result indicates that *M. genomsp.* were more closely related than expected by chance across women who delivered preterm. It suggests that sPTB-associated vaginal conditions across women may be more conserved, harboring a group of significantly closely related *M. genomsp.* strains, compared to TB-associated vaginal conditions.

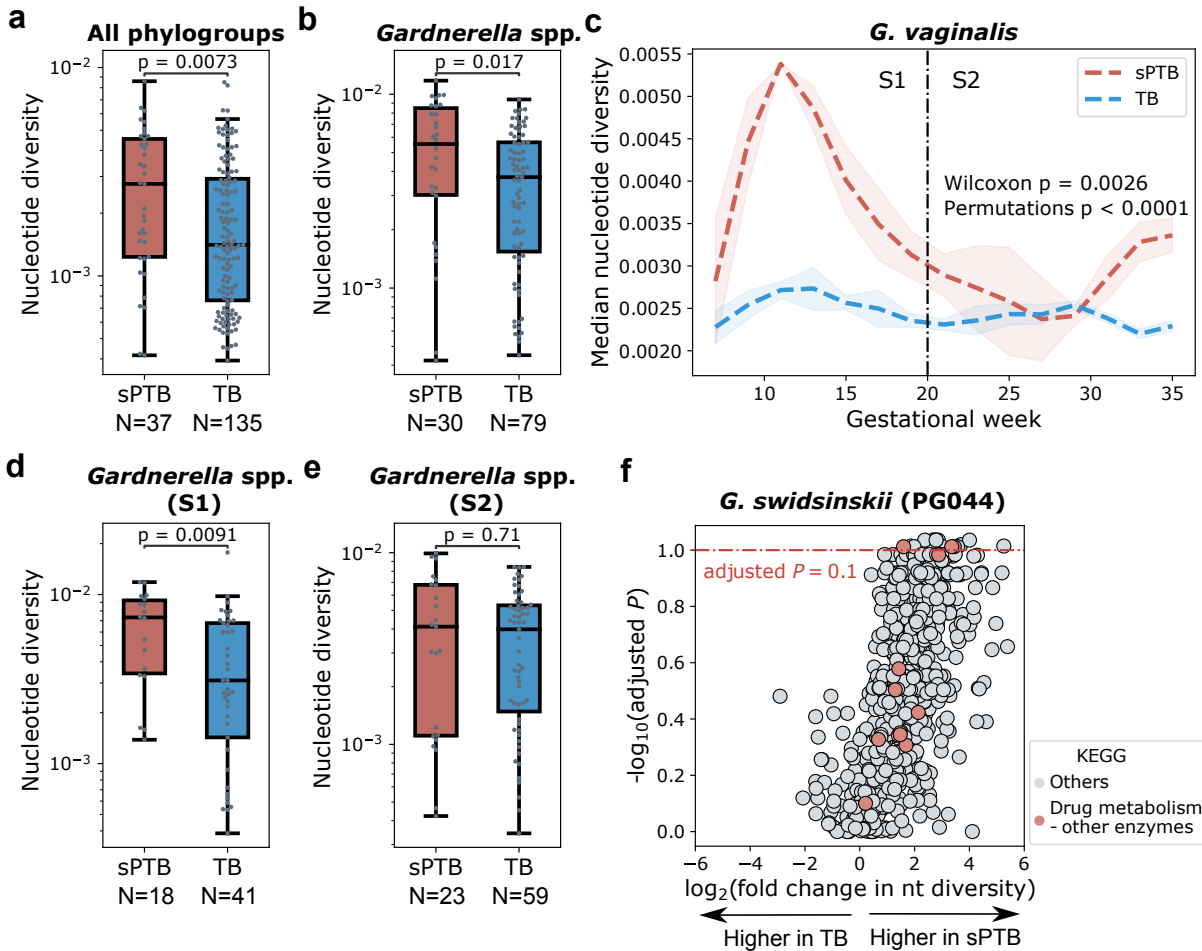


**Fig. 1 | The composition of the vaginal microbiome associated with sPTB. a.** Phylogenetic tree of non-redundant MAGs representing 132 species-level phylogroups differing by at least 95% average nucleotide identity (ANI) based on concatenated amino acid (AA) sequences of 120 marker genes. Representative MAGs of 25 phylogroups had <60% of marker genes AA sequence identified and were not included in the tree. Gray nodes indicate a bootstrap value >80. The tree is rooted by midpoint and annotated by the GC content and genome size of the representative MAGs. **b.** Compositional tensor factorization (CTF) analysis showing microbiome composition trajectories over gestational ages separated by pregnancy outcomes using the top two ordination axes (Component 1 and Component 2). Each dot represents a subject. **c.** Component 1 in the CTF analysis compared between sPTB and TB. Box, IQR; line, median; whiskers, 1.5\*IQR;  $p$ , two-sided Mann-Whitney. **d.** Feature rankings of phylogroups colored by preterm birth (sPTB) and term birth (TB) based on Component 1 in the CTF analysis

## ***Gardnerella* species have higher microdiversity in the first half of pregnancies that ended preterm.**

Human microbes can adapt to host-induced environmental changes (e.g., diet, antibiotics) through genetic variations<sup>29</sup>. Therefore, the microbial populations of the same species in different hosts can have a different genetic structure, which provides them a competitive advantage. These genetic differences, in turn, may be related to the phenotype of the host. To understand the genetic structure of microbial populations in the vaginal environment and its association with pregnancy outcomes, we calculated the nucleotide diversity for each identified phylogroup. Overall, vaginal microbial populations had a significantly higher genome-wide nucleotide diversity in sPTB than in TB (median along pregnancy; Mann-Whitney  $P = 0.0073$ ; **Fig. 2a**). Stratifying by phylogroups, we found that this difference was mainly driven by *Gardnerella* spp. ( $P = 0.017$ ; **Fig. 2b**). *G. piotti* (PG042), *G. swidsinskii* (PG044), and *G. gs13* (PG051) and a potentially novel *Gardnerella* spp. (PG045), along with a phylogroup of *Atopobium vaginae*, a suspected vaginal pathobiont<sup>30</sup>, showed significantly higher genome-wide nucleotide diversity in sPTB ( $P < 0.05$ , adjusted  $P < 0.1$  for all; **Supplementary Fig. 3**). These results imply that microbial populations composed of more diverse strains from the same species, and particularly *Gardnerella* spp., are growing in the vaginal environment associated with sPTB.

To understand how the nucleotide diversity of *Gardnerella* spp. changes over time during pregnancy, we analyzed temporal trajectories of term and preterm pregnancies. To this end, we pooled the data from all women in each group, binned pregnancy weeks and used splines to smooth the temporal curves (Methods). We found a significant difference between the temporal trajectories of *Gardnerella* spp. nucleotide diversity in pregnancies ending at term and preterm (Permutation test  $< 0.001$  (ref. <sup>31</sup>), Wilcoxon signed-rank  $P < 0.003$ ; Methods; **Fig. 2c**). Specifically, we found that the nucleotide diversity of *Gardnerella* spp. increased at the beginning of pregnancies which ended preterm, with a peak at around gestational week 13, and then dropped to its initial value at around gestational week 20 (**Fig. 2c**). In comparison, nucleotide diversity of *Gardnerella* spp. in TB remained relatively stable (**Fig. 2c**). Given that gestational week 20 is the middle of a full-term pregnancy, we subsequently analyzed samples with respect to two time periods - first half (0-19 gestational week) and second half of pregnancy (20-37 gestational week; 37 was chosen to ensure a similar time range for both sPTB and TB). As expected, the nucleotide diversity of *Gardnerella* spp. in sPTB was significantly higher than TB in the first half of pregnancy (median along first half; Mann-Whitney U  $P = 0.0091$ ; **Fig. 2d**), but not in the second half ( $P = 0.71$ ; **Fig. 2e**). We further found that nucleotide diversity had a significantly stronger correlation with synonymous mutations than with nonsynonymous mutations across *Gardnerella* spp. (paired t-test  $P = 0.0011$ ; **Supplementary Fig. 4**), suggesting a more important role of purifying selection in shaping genomic diversity. Overall, these results suggest that genetic diversity of *Gardnerella* spp. in the first half of pregnancy is important to birth outcomes, and could perhaps be used as a biomarker for early diagnosis of sPTB.



**Fig. 2 | Microdiversity patterns of the vaginal microbiome are associated with sPTB.** **a-b**, A comparison of median genome-wide nucleotide diversity along pregnancy between sPTB and TB, displayed for all phylogroups (a) and *Gardnerella* spp. (b). **c**, Trajectory of median nucleotide diversity of *Gardnerella* spp. along pregnancy. S1 - first half of pregnancy; S2 - second half of pregnancy. The shaded area depicts mean  $\pm$  s.d./n. **d-e**, A comparison of median genome-wide nucleotide diversity between sPTB and TB, displayed for pregnancy S1 (d) and S2 (e). **f**, Volcano plot illustrating the significance (Mann-Whitney; y-axis) of difference between nucleotide diversity (fold change; x-axis) in sPTB and TB of every gene in *G. swidsinskii* (PG044). Genes above the red dashed line have  $P < 0.05$  and an adjusted  $P < 0.1$ . Genes belonging to KEGG pathways that were significantly enriched in genes showing significant nucleotide diversity differences are color-coded (adjusted  $P < 0.1$ ).

To understand if any particular genes are driving the association between sPTB and the microdiversity of *Gardnerella* spp. in the first half of pregnancy, we further analyzed nucleotide diversity at the gene level for these species. We identified 21 and 47 genes (out of 825 and 531) in *G. swidsinskii* (PG044) and *G. vaginalis* (PG043), respectively, that showed significantly different nucleotide diversity between sPTB and TB (median along the first half of pregnancy; Mann-Whitney  $P < 0.05$ , adjusted  $P < 0.1$  for all). These genes included one gene encoding the putative tail-component of bacteriophage HK97-gp10 ( $P = 0.0012$ ) and one gene encoding

putative AbiEii toxin, Type IV toxin–antitoxin system ( $P = 5 \times 10^{-4}$ ), which might be involved in the interaction with maternal health<sup>32,33</sup>. To further identify what functions were related to these associations, we then performed functional enrichment analysis (Methods) using the eggNOG functional annotation of genes (**Supplementary Table 2**). We found that the KEGG pathway ‘drug metabolism - other enzymes’ (ko00983) was significantly enriched among genes from *G. swidsinskii* (PG044) that had significantly higher microdiversity ( $P < 0.05$ , adjusted  $P < 0.1$ ; **Fig. 2f**). This result suggests that the more diverse gene pool in *G. swidsinskii* (PG044) detected in sPTB may be associated with adaptation to drugs present in the environment. This may be consistent with our recent finding that xenobiotics detected in the vaginal environment are strongly associated with sPTB<sup>34</sup>.

To verify that the higher nucleotide diversity we observed in sPTB pregnancies was not caused by sampling or sequencing bias, we compared the read count and quality of MAGs obtained from sPTB and TB samples. If this higher diversity is the result of a higher read count in sPTB samples or more complete MAGs, we would expect read count and MAG completeness to be higher in sPTB samples. Instead, we found the completeness and contamination of MAGs assembled from sPTB samples were not significantly different from TB (Mann-Whitney  $P = 0.71$  and  $0.73$ , respectively; **Supplementary Fig. 5a,b**). Next, we assessed the correlation between the number of reads mapped to each phylogroup and its genome-wide diversity. If a higher diversity is caused by more reads mapped to the MAG representing the phylogroup, we would expect a positive correlation between these two measurements. However, only 3 phylogroups (PG064, a *Dialister* spp.; PG102, a *Peptoniphilus* spp.; and PG122, a *Bradyrhizobium* spp.) had a significant positive correlation between read counts and nucleotide diversity (Spearman  $\rho = 0.70, 0.54, \text{ and } 0.42$ , respectively; non-adjusted  $P = 0.035, 0.024, \text{ and } 0.00015$ , respectively). In 98% of phylogroups, we did not observe a statistically significant positive correlation (median [IQR] Spearman correlation of  $-0.067 [-0.19, -0.13]$ ). None of the four *Gardnerella* spp. Phylogroups that showed significantly higher nucleotide diversity in sPTB pregnancies in **Supplementary Fig. 3** were significantly positively correlated (Spearman  $\rho = -0.00, -0.12, -0.02, \text{ and } -0.35$  and  $P = 0.96, 0.091, 0.77, \text{ and } 0.00045$  for PG042, PG044, PG045, and PG051, respectively; **Supplementary Fig. 5c**).

Finally, as we have observed a significantly higher read count in sPTB samples (Mann-Whitney  $P = 0.0004$  and  $0.061$  for samples and subjects, respectively; **Supplementary Fig. 5d and 5e**, respectively), we subsampled an identical number of reads ( $10^5$ ) from each sample, retaining 75% of samples, and repeated our analyses of nucleotide diversity. As with the first analysis (**Fig. 2**), nucleotide diversity was significantly higher in sPTB pregnancies across all phylogroups, and particularly in *Gardnerella* spp. (Mann-Whitney  $P = 0.015$  and  $P = 0.0043$ , respectively; **Supplementary Fig. 5f and 5g**, respectively). Similarly, in the first half of pregnancy, the nucleotide diversity of *Gardnerella* spp. was significantly higher in sPTB ( $P = 0.026$ ; **Supplementary Fig. 5h**), while in the second half, there was no significant difference ( $P = 0.22$ ; **Supplementary Fig. 5i**). Overall, these results confirm that the sPTB-associated nucleotide diversity we observed was not biased by technical artifacts.



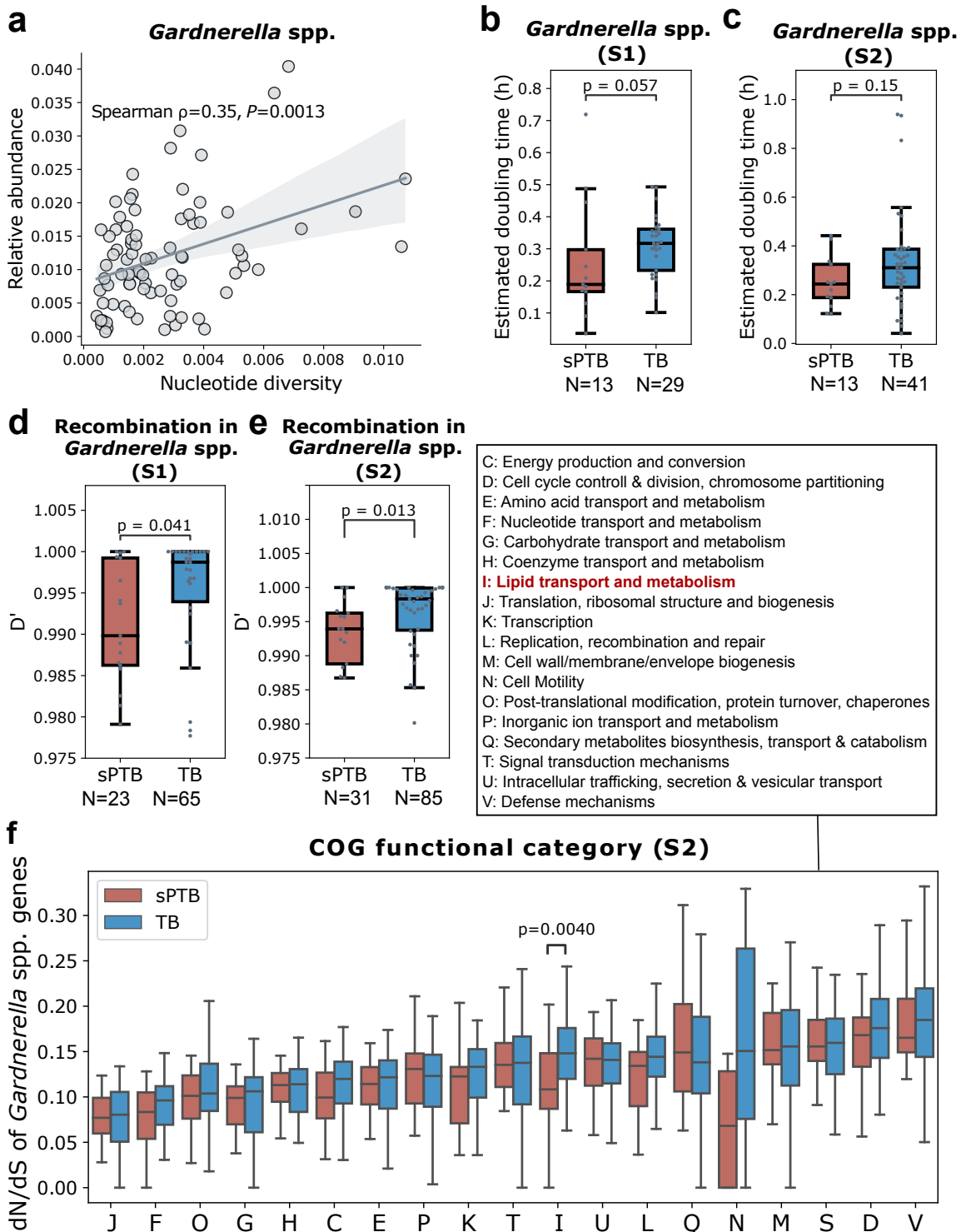
## Evolutionary forces acting on *Gardnerella* species are associated with pregnancy outcomes.

Adaptation should increase the fitness of an organism, its ability to survive and reproduce in a given environment. To better understand if the *Gardnerella* spp. populations with higher genetic diversity grow better in the vaginal environment associated with sPTB, we inferred fitness using two measures: relative abundance and growth rate. Indeed, we found that nucleotide diversity in these species was positively correlated with relative abundance (Spearman  $\rho = 0.35$ ,  $P = 0.0013$ ; **Fig. 3a**). This correlation was not observed in other phylogroups (**Supplementary Fig. 6a**). *L. crispatus* (PG080) and *L. iners* (PG086) even showed a significantly negative correlation ( $\rho = -0.39$  and  $-0.32$ ,  $P = 0.026$  and  $0.0014$ , respectively; **Supplementary Fig. 6b,c**). We additionally used gRodon<sup>35</sup> to predict the maximum growth rate of microbes based on codon usage bias in highly expressed genes encoding ribosomal proteins. We found that in the first half of pregnancy, *Gardnerella* spp. had a somewhat higher, albeit not statistically significant, maximum growth rate in sPTB pregnancies (Mann-Whitney  $P = 0.057$ ; **Fig. 3b**), while in the second half of pregnancy, the difference was diminished ( $P = 0.15$ ; **Fig. 3c**). These results suggest that the sPTB-associated genetic diversity observed in *Gardnerella* spp. may be related to the optimization for faster growth in the sPTB-associated vaginal environment.

Microbial population structure is influenced by various evolutionary processes including selection and homologous recombination<sup>36</sup>. Competence, a mechanism of horizontal gene transfer which involves homologous recombination, has been identified in *Gardnerella* spp.<sup>37</sup>. To better interpret the significant differences we observed in the microdiversity patterns of *Gardnerella* spp. between sPTB and TB, we quantified the degree of homologous recombination using the normalized coefficient of linkage disequilibrium between alleles at two loci,  $D'$ . A value of  $D'$  closer to 0 indicates a higher degree of recombination<sup>38</sup>. Interestingly, we found that the median  $D'$  of *Gardnerella* spp. was significantly smaller in sPTB pregnancies in both the first (Mann-Whitney  $P = 0.041$ ; **Fig. 3d**) and second halves of pregnancy ( $P = 0.013$ ; **Fig. 3e**), and the same was also observed for the  $D'$  of three specific *Gardnerella* spp., *G. piotti* (PG042), *G. gs7* (PG046), and PG047, in the first half of pregnancy ( $P < 0.05$ , adjusted  $P < 0.1$  for all; **Supplementary Fig. 7a**). This result suggests that *Gardnerella* spp. tends to have more frequent recombination in women who delivered preterm during both halves of pregnancy.

Next, we quantified the degree of selection using dN/dS in this species (Methods). This measure quantifies the ratio between synonymous and non-synonymous mutations, and hence offers insight into the type of selection, with values close to zero indicating purifying selection, and values higher than one indicating positive selection<sup>39</sup>. dN/dS is calculated in relation to the reference, and can therefore detect selection on mutations that have already been fixed within the population<sup>40</sup>. Consistent with the gut and ocean microbiomes<sup>41–43</sup>, purifying selection is predominant across all genes of the vaginal microbiome (dN/dS  $\ll 1$ ; median [IQR] dN/dS of 0.17 [0.10, 0.29]; **Supplementary Fig. 7b**). While the median dN/dS of all *Gardnerella* spp. genes was not significantly different between sPTB and TB pregnancies (Mann-Whitney U  $P = 0.48$ ), we detected some differences when examining high-level functions (COG categories<sup>44</sup>) within each half of pregnancy. In the first half, the median dN/dS of *Gardnerella* spp. genes was

somewhat lower in sPTB pregnancies for inorganic ion transport and metabolism, lipid transport and metabolism, secondary structure, and cell wall/membrane/envelope biogenesis, though this was not statistically significant after adjusting for multiple testing (COG categories “P”, “I”, “Q”, and “M”, respectively; Mann-Whitney  $P < 0.05$ , adjusted  $P > 0.1$  for all; **Supplementary Fig. 7c**). In the second half of pregnancy, the median dN/dS was significantly lower in sPTB pregnancies for lipid transport and metabolism and cell motility (COG categories “I” “N”; Mann-Whitney  $P = 0.0040$  and  $P = 0.04$ , adjusted  $P = 0.07$  and  $0.40$ , respectively; **Fig. 3f**). Our results suggest that *Gardnerella* spp. genes involved in lipid transport and metabolism may undergo stronger purifying selection in sPTB. As purifying selection maintains the fitness of organisms by constantly sweeping away deleterious mutations and conserving functions, *Gardnerella* spp. may benefit from this stronger purifying selection targeting lipid functioning when growing in the sPTB-associated vaginal environment during pregnancy.



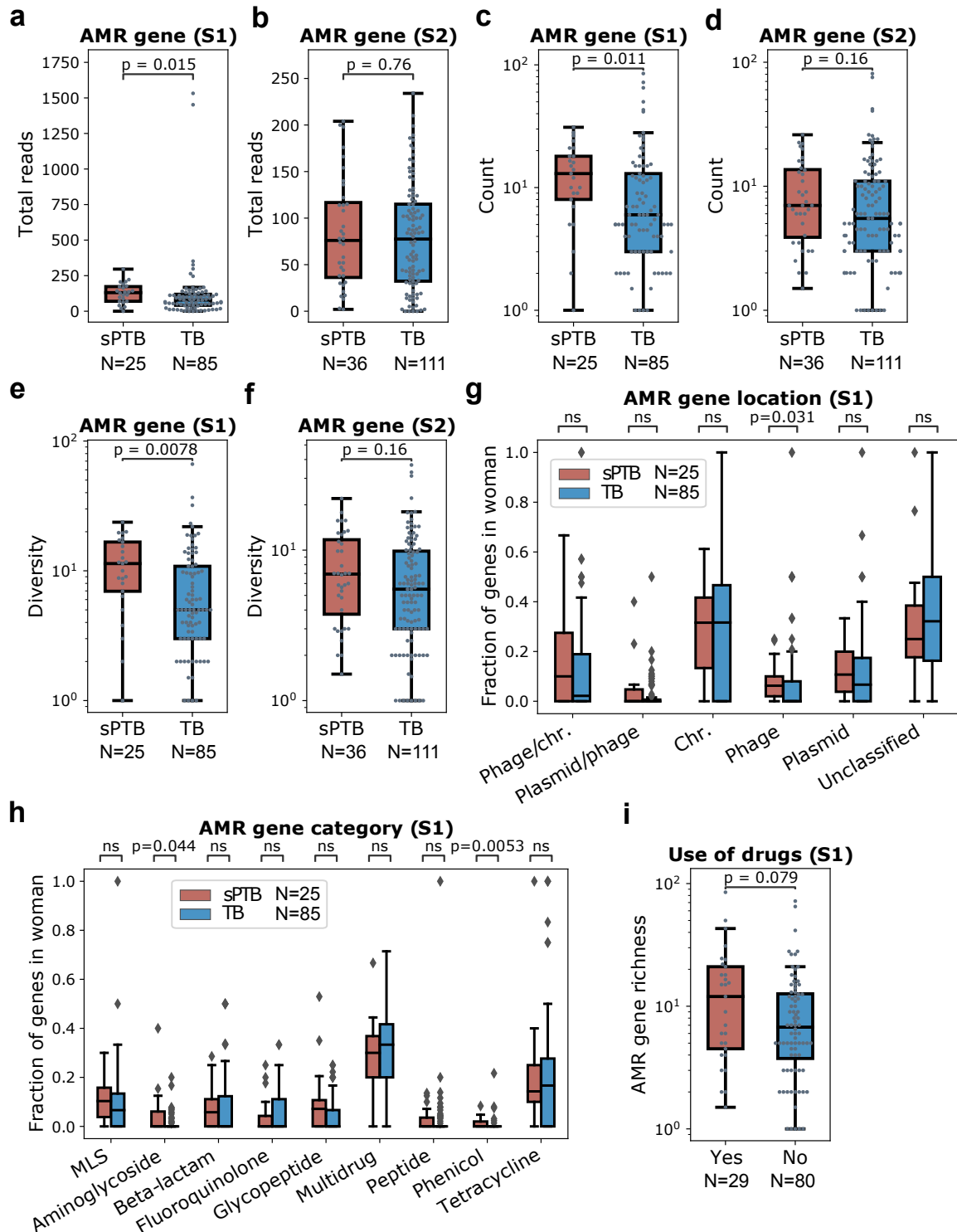
**Fig. 3 | Evolutionary forces on *Gardnerella* spp.** **a.** Spearman correlation between median genome-wide nucleotide diversity and relative abundance of *Gardnerella* spp. along pregnancy. The line and the shaded area depict the best-fit trendline and the 95% confidence interval (mean  $\pm$  1.96 s.e.m.) of the linear regression. **b,c,** Predicted maximal doubling time (gRodon<sup>35</sup>) of *Gardnerella* spp. compared

between sPTB and TB, displayed for the first (b, S1) and second (c, S2) halves of pregnancy. Box, IQR; line, median; whiskers, 1.5\*IQR;  $p$ , two-sided Mann-Whitney. **d,e**, Median  $D'$  of *Gardnerella* spp compared between sPTB and TB, displayed for the first (S1, **d**) and second (S2, **e**) halves of pregnancy. Lower  $D'$  indicates more frequent recombination. **f**, dN/dS of *Gardnerella* spp genes compared between sPTB and TB by COG functional categories, displayed for the second half of pregnancy (S2). dN/dS closer to 0 indicates stronger purifying selection.

### **sPTB-associated vaginal microbiomes have a higher antibiotic-resistance potential.**

Antibiotics are widely used during pregnancy, sometimes even topically in the vagina<sup>45</sup>. This exposure may promote antimicrobial resistance (AMR). To assess if antibiotic-resistance potential in the vaginal microbiome is associated with sPTB, we subsampled an identical number of reads ( $10^5$ ) from each sample and mapped them to the Comprehensive Antibiotic Resistance Database<sup>46</sup>. The total number of reads mapped to AMR reference genes was significantly higher in the first half of sPTB pregnancies (Mann-Whitney U  $P = 0.015$ ; **Fig. 4a**), but not in the second half ( $P = 0.76$ ; **Fig. 4b**). In addition, to assess the difference of specific AMR genes between the vaginal microbiomes of sPTB and TB, we identified AMR genes in the genomic assemblies. A significantly higher median count and Shannon-Wiener diversity of AMR genes were detected in vaginal microbes sampled at the first half of pregnancies that ended preterm (3-times higher on average; Mann-Whitney U  $P = 0.011$  and  $P = 0.0078$ , respectively; ; **Fig. 4c and 4e**, respectively), yet this difference was not detected in the second half ( $P = 0.16$  for both; **Fig. 4d and 4f**, respectively). Exploring the source of these genes, we found a significantly higher median fraction of phage-borne AMR genes in the microbiomes of women who delivered preterm ( $P = 0.031$ ; **Fig. 4g**), suggesting transduction may promote the higher median count and diversity of AMR genes observed in the first half of sPTB pregnancies (**Fig. 4c,e**). Among the 9 AMR gene categories that had genes present in at least 10% of women, phenicol and aminoglycoside resistance genes showed a significantly higher median fraction in the sPTB microbiome ( $P = 0.041$  and  $P = 0.032$ , respectively; **Fig. 4h**). These results suggest a unique antibiotic resistance profile associated with the first half of sPTB pregnancies, potentially indicative of usage of specific antibiotics. Indeed, we detected a somewhat higher richness of AMR genes along the first half of pregnancy in women who used antibiotics in the past 6 months before pregnancy than those who did not (Mann-Whitney U  $P = 0.079$ ; **Fig. 4i**). This is also consistent with our observation that genes with sPTB-associated nucleotide diversity were enriched for drug metabolism in *G. swidsinskii* (PG044) (**Fig. 2d**).

To check if the strong association between sPTB and the AMR potential of the vaginal microbiome is contributed by a particular phylogroup, we performed a similar analysis for each phylogroup. We found, however, that none of them showed a significant difference in the median count and diversity of AMR genes between sPTB and TB (Mann-Whitney U  $P > 0.05$  for all). This result suggests that the higher AMR potential associated with sPTB may be a property of the vaginal microbiome as an ecosystem. However, this lack of association could also be driven by underestimation of AMR genes due to the limitation of MAG binning methods in recovering mobile genetic elements<sup>47</sup>.



**Fig. 4 | Antimicrobial resistance (AMR) gene profiles of the vaginal microbiome are associated with sPTB. a,b,** Total subsampled reads ( $10^5$ ) mapped to AMR genes compared between sPTB and TB, in the first (S1, **a**) and second (S2, **b**) halves of pregnancy. **c,d,** Median count (along period) of AMR genes compared between sPTB and TB, in the first (S1, **c**) and second (S2, **d**) halves of pregnancy. **e,f,** Median Shannon-Wiener diversity (along period) of AMR genes compared between sPTB and TB, in the first (S1,

**e**) and second (S2, **f**) halves of pregnancy. **g**, Fraction of AMR genes originating in different locations, shown as median along the first half of each pregnancy. Chr.: chromosome. **h**. Fraction of AMR genes belonging to different resistance categories, shown as median along the first half of each pregnancy. MLS, macrolide, lincosamide and streptogramin B. **i**. AMR gene richness in the first half of pregnancy (S1) compared between women who used and did not use drugs in the past 6 months before pregnancy. Box, IQR; line, median; whiskers, 1.5\*IQR; *p*, two-sided Mann-Whitney U.

## Discussion

Microbial genomes can exhibit large variations even within the same species, as a result of adaptation to various environments<sup>48</sup>. Associations between the vaginal microbiome and preterm birth have been widely reported<sup>7,8,12,49</sup>. However, there is still much left to explore regarding potential mechanisms underlying host-microbiome interactions in this context. Here, by leveraging publicly available metagenomic data<sup>1</sup>, we provide a population genetic view of the vaginal microbiome during pregnancy. We identify a number of novel microbial features including population nucleotide diversity, selection metrics, and antibiotic resistance potential that are associated with sPTB. Interestingly, we find that the higher population nucleotide diversity is driven by *Gardnerella* spp. during the first half of pregnancy. This species appears to undergo more intense changes in the population structure contributed by recombination and purifying selection in pregnancies which ended preterm. We also show evidence that this sPTB-associated genetic pattern of *Gardnerella* spp. may be related to optimization of growth rates in vaginal conditions linked to sPTB. Our results are indicative of adaptation of the vaginal microbiota to the host, which in turn may influence pregnancy outcomes.

Our findings regarding a relationship between ecological processes in the pregnancy vaginal microbiome and subsequent preterm birth are consistent with previous studies<sup>1,7,11,12,14,50,51</sup>. We add to these previous studies by exploring an additional layer of microbial variability associated with sPTB - microbial genetic diversity. It is known that genomic variation within species can result in phenotypic diversity and adaptations to different environments<sup>48</sup>. These adaptations, in turn, can affect host phenotypes such as disease outcomes<sup>52</sup>. Such associations between microbial genomic variation and host phenotypes have been reported in the gut microbiome<sup>42,53,54</sup>. Our study suggests that this phenomenon also occurs in the vaginal ecosystem, and suggests that it may be associated with pregnancy outcomes. Nevertheless, the associations between microbial genetic diversity and pregnancy outcomes we detect might also be a consequence of a different process that acts on both variables, and while we find this unlikely, this should be determined by future studies.

Interestingly, we found that the association of genetic diversity and sPTB was largely driven by *Gardnerella* spp., a group of species commonly associated with BV<sup>50,55,56</sup>. A number of studies reported a higher abundance of these species in sPTB pregnancies<sup>9,50,57,58, 1,9</sup>. We show that *Gardnerella* spp. populations with more genetically diverse strains may also be associated with sPTB. In addition, we found that this taxon has the capacity to grow 1.5 times faster in

pregnancies that ended preterm, consistent with an overall higher relative transcriptional rate of *G. vaginalis* which was previously reported<sup>1</sup>. These more genetically diverse strains appear to have adapted to the vaginal environment associated with sPTB, exhibiting higher fitness. Notably, the higher genetic diversity associated with sPTB in *Gardnerella* spp. was detected during the first half of the pregnancy (<20 gestational week) rather than the second half. Most potential biomarkers of sPTB (e.g., serum alpha-fetoprotein<sup>59</sup> were so far identified using samples from the second trimester of pregnancy (gestational week 14-27). Our results suggest that high resolution analysis of microbiome samples from even earlier stages of pregnancy (<week 20) may yield informative biomarkers of pregnancy outcomes.

As in the human gut microbiome<sup>19-21</sup>, we show evidence that adaptive evolution also occurs in the vaginal microbiome. Several environmental factors affecting the vaginal ecosystem, such as pH, neutrophil levels, and xenobiotics, have been reported to be associated with sPTB<sup>34,60</sup>. These environmental factors may act as selective stressors that lead to different evolutionary patterns in the vaginal microbiome. Indeed, we detected more frequent homologous recombination and stronger purifying selection within *Gardnerella* spp. during pregnancies that end preterm. Homologous recombination is a critical mechanism speeding adaptation by increasing fixation probability of beneficial mutations<sup>61</sup> and reducing clonal interference (i.e., competition between beneficial mutations) in bacteria<sup>62</sup>. Purifying selection also contributes to adaptation by sweeping away deleterious mutations and conserving functions, such as in oligotrophic nutrient conditions<sup>41,63</sup>. Notably, we found that sPTB-associated purifying selection is particularly strong on genes involved in lipid transportation and metabolism. This is consistent with previous identification of lipid metabolites (e.g., monoacylglycerols and sphingolipids) as signatures of sPTB<sup>34,64,65</sup>. Whether this stronger purifying selection targeting lipid transportation and metabolisms in pregnancy that ended preterm leads to changes in the concentrations of lipid metabolites however requires further experimental testing. As both recombination and purifying selection can reduce genetic diversity, sPTB-associated recombination and purifying selection along pregnancy may explain the higher nucleotide diversity of *Gardnerella* spp. in sPTB in the first half of pregnancy compared to the second half.

Antibiotics are common selective stresses acting on the human microbiome<sup>29</sup> and have been associated with preterm birth<sup>66</sup>. We detected higher count and diversity of AMR genes associated with sPTB, which our analysis suggests to be facilitated by prophages in preterm vaginal microbiomes. While multiple phages (e.g., Siphoviridae, Myoviridae, and Microviridae) have been detected in the vagina of pregnant women, their association with sPTB is rarely studied<sup>67</sup>. Our results imply a potentially important role of bacteria-phage interactions in pregnancy outcomes via transferring of AMR genes. We also found that genes related to phenicol and aminoglycoside resistance were more abundant in vaginal microbiomes during pregnancies that ended preterm. While both antibiotics have been frequently used to treat gynecologic infection for decades<sup>68</sup>, and some phenicols (e.g., chloramphenicol) are thought to be safe for use during pregnancy<sup>69</sup> aminoglycoside is teratogenic. Previous studies reported that exposure to antibiotics could change the composition of the vaginal microbiome<sup>45,70</sup>, indicating an ecological effect. In comparison, our results may suggest adaptation of the vaginal microbiome to more frequent antibiotics usage in women who delivered preterm, leading to an

enrichment of AMR genes as well as higher nucleotide diversity in *Gardnerella* spp. genes encoding enzymes for drug metabolism. While this hypothesis requires further study, it is further supported by the fact that a higher proportion of women who delivered preterm (31%) had used antibiotics in the past 6 months before pregnancy than women who delivered at term (23%).

Despite its findings, our study is limited by low sequencing depth (median bacterial read count <math>5 \times 10^5</math>) and inconsistent sampling frequency during pregnancy (1 to 8 samples per pregnancy, with an average of 4). These limitations lead to high sparsity in the features analyzed, preventing a more in-depth temporal and predictive analysis of the link between population genetics of the vaginal microbiome and sPTB. Our results warrant a high-resolution investigation of the vaginal metagenome, with frequent sampling and high sequencing depth.

In summary, through in-depth population genomic analyses, our study identified novel genetic and functional associations between the vaginal microbiome and preterm birth. We revealed evidence of microbial genetic adaptation to the host environment linked to preterm birth and highlighted the importance of microbial evolutionary processes to adverse pregnancy outcomes, particularly in *Gardnerella* spp.. Future investigation on the pressures driving the sPTB-associated microbial adaptation is warranted to fully understand the molecular mechanisms underlying preterm birth.

## Methods

### Sample selection and metagenomic data

Metagenomic sequencing data<sup>1</sup> generated from 135 vaginal samples collected longitudinally during pregnancy from 40 women with majority of them identifying as Black women who eventually delivered preterm spontaneously (sPTB) and 570 vaginal samples from 135 women who delivered at term (TB) were obtained from dbGaP (study no. 20280; accession ID phs001523.v1.p1). We used the same definitions for preterm birth as in Fettweis et al. 2019<sup>1</sup>: spontaneous preterm birth is defined as live birth between 23 and 37 gestational weeks without medical indication, and term birth is defined as live birth at or after 39 gestational weeks. To check for the presence of some potential confounders for vaginal microbiome-sPTB associations, we calculated propensity scores<sup>71</sup> for each subject based on income, age, and race using a logistic regression model. We found that propensity scores for both sPTB and TB subjects exhibited a similar distribution (Kolmogorov–Smirnov test  $P = 0.21$ ), suggesting the associations we detect are not likely to be confounded with these variables (**Supplementary Fig. 8**). These results suggest a negligible confounding effect of income, age, and race in this study on microbiome-sPTB associations.

### Metagenomic assembly, genomic binning, genome annotation, and relative abundance

Our analysis follows the accepted standards used in refs<sup>72–75</sup>, using the ATLAS pipeline<sup>76</sup>. Bases with quality scores <math>< 25</math>, raw reads <math>< 50</math> bp lengths, and sequencing adapters were removed using Trimmomatic v.0.39<sup>77</sup>. Reads mapped to human and PhiX genome sequences were removed by mapping with Bowtie2 v.2.3.5.1<sup>78</sup>. Assembly and binning were done with ATLAS: filtered reads were assembled using metaSPAdes v.3.15.2<sup>79</sup>, and contigs were binned



into metagenome-assembled genomes (MAGs) using MetaBAT2 v.2.14.0<sup>80</sup> with a minimum contig length of 1500. Quality, GC content, genome size, and taxonomy of MAGs were estimated using CheckM v.1.0.9<sup>24</sup>. MAGs were de-replicated using dRep v.3.2.0<sup>81</sup> with an average nucleotide identity (ANI) of 0.95, minimum completeness of 50%, and maximum genome contamination of 10%. The MAG with the highest dRep score within each 95% ANI cluster, termed here as a phylogroup, was selected as the representative MAGs for that phylogroup. Genes were predicted using Prodigal v.2.6.3<sup>82</sup> and annotated using EggNOG v.5.0<sup>83</sup>. Filtered reads were mapped to representative MAGs using Bowtie2 v.2.3.5.1<sup>84</sup>. The relative abundance of each representative MAG was calculated by dividing the number of reads that mapped to that MAG, corrected to the genome size and completeness, by the total number of reads in each sample.

### Tensor factorization

To characterize and compare the dynamics of the vaginal microbiome in term and preterm pregnancies, we used a revised version of compositional tensor factorization<sup>27</sup>

### Phylogeny

Amino acid (AA) sequences of 120 marker genes were called and aligned for representative MAGs using GTDB-Tk v.1.5.1<sup>85</sup>. MAGs with <60% of AA in the alignment were excluded in the phylogenetic tree construction. The best evolutionary model LG+G+I (the Le Gascuel model + gamma distribution + invariant sites) was identified using protest3 v.3.4.2<sup>86</sup> and 500 bootstraps were used for tree construction using RAxML v.8.2.12<sup>87</sup>. The tree was rooted by midpoint and visualized in iTol v.6.3<sup>88</sup>.

### Microdiversity profiling, growth rate estimation, and antimicrobial resistance genes

Population microdiversity metrics including genome-wide nucleotide diversity, gene-wide nucleotide diversity, linkage disequilibrium measures ( $D'$ ) and dN/dS, were calculated using InStrain v1.0.0<sup>40</sup> using the 157 representative MAGs as the reference database. Maximum growth rate was estimated for each MAG using gRodon<sup>35</sup>. Antimicrobial resistance (AMR) genes were detected in assemblies and MAGs using PathoFact v.1.0<sup>89</sup> with default parameters.

### Functional enrichment analysis

To identify COG/KEGG pathways that were enriched in genes showing significant difference in nucleotide diversity between sPTB and TB, the frequency of each COG/KEGG category was first calculated from significant genes (observed frequency). Then, the frequency of each COG/KEGG category was calculated from an identical number of genes randomly selected from all genes (expected frequency). This process was repeated 10,000 times. The null hypothesis was that the observed frequency of COG category is smaller than the expectation. For each COG, probability  $P$  of the null hypothesis was calculated using the formula:  $P = |\{x_i \in x : x_i \geq k\}| / 10000$ , where [...] denotes a multiset,  $x = (x_1, x_2, \dots, x_n)$  is a list of expected values, and  $k$  is the observed value.

### Statistical analysis

A different number of samples was available for each woman in the database. In our analyses, we therefore used the median along pregnancy (or its first or second half). The false-discovery

rate procedure (FDR) of Benjamini and Hochberg (BH)<sup>90</sup> was used to correct for multiple testing. Adjusted  $P < 0.1$  was used as the significance cutoff.

### Temporal analysis

To generate the trajectories representing the change in nucleotide diversity over time in term and preterm deliveries, we pooled the temporal data of *Gardnerella* spp. from all women in each group (term and preterm). When we had more than one observation per gestational week, we took the median value across samples. We then binned the temporal data into bins of 3 weeks, except for the first bin which spanned weeks 1-7, and took the median of each bin as a summary. To smooth the observed binned data we applied splines, which is a special function defined piecewise by polynomials for data smoothing. To compare between the temporal trajectories of preterm and term, we performed a permutation test, in which we generated a null distribution of euclidean distances by shuffling these trajectories  $10^4$  times and comparing to the euclidean distance in the original data<sup>31</sup>.

### Data availability

The dataset used is available from dbGaP (phs001523).

### Acknowledgements

We thank members of the Korem lab for useful discussions. This study was supported by the Eunice Kennedy Shriver National Institute of Child Health and Human Development (NICHD) of the National Institutes of Health under award number R01HD106017, the Program for Mathematical Genomics at Columbia University, and the CIFAR Azrieli Global Scholarship in the Humans & the Microbiome Program (T.K.). The dataset used was obtained from dbGaP (phs001523), using data provided by Gregory A. Buck, Ph.D. and colleagues and supported by NICHD (U54 HD080784) (G.A.B).

### Author contributions

J.L. and T.K. designed the study. J.L. and L.S. analyzed the data with input from T.K., M.S., and G.A.B. J.L. wrote the manuscript with input from L.S., T.K., M.S., B.Z. and G.A.B. G.A.B. assisted with data access and acquisition.

### Competing interests

G.A.B. is a member of the Scientific Advisory Board of Juno, LTD., a startup biotech firm focused on using the vaginal microbiome to address issues of women's gynecologic and reproductive health. Juno had no involvement in the current study. Other authors declare no competing interests.

### References

1. Fettweis, J. M. *et al.* The vaginal microbiome and preterm birth. *Nat. Med.* **25**, 1012–1021 (2019).

2. Tiensuu, H. *et al.* Risk of spontaneous preterm birth and fetal growth associates with fetal SLIT2. *PLoS Genet.* **15**, e1008107 (2019).
3. Preterm birth. <https://www.who.int/news-room/fact-sheets/detail/preterm-birth>.
4. Goldenberg, R. L., Culhane, J. F., Iams, J. D. & Romero, R. Epidemiology and causes of preterm birth. *Lancet* **371**, 75–84 (2008).
5. Hong, X. *et al.* Genome-wide approach identifies a novel gene-maternal pre-pregnancy BMI interaction on preterm birth. *Nat. Commun.* **8**, 15608 (2017).
6. Hong, X. *et al.* Genome-wide association study identifies a novel maternal gene $\times$  stress interaction associated with spontaneous preterm birth. *Pediatr. Res.* 1–8 (2020).
7. DiGiulio, D. B. *et al.* Temporal and spatial variation of the human microbiota during pregnancy. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 11060–11065 (2015).
8. Ravel, J. *et al.* Vaginal microbiome of reproductive-age women. *Proc. Natl. Acad. Sci. U. S. A.* **108 Suppl 1**, 4680–4687 (2011).
9. Tabatabaei, N. *et al.* Vaginal microbiome in early pregnancy and subsequent risk of spontaneous preterm birth: a case-control study. *BJOG* **126**, 349–358 (2019).
10. Chu, D. M., Seferovic, M., Pace, R. M. & Aagaard, K. M. The microbiome in preterm birth. *Best Pract. Res. Clin. Obstet. Gynaecol.* **52**, 103–113 (2018).
11. Freitas, A. C., Bocking, A., Hill, J. E., Money, D. M. & VOGUE Research Group. Increased richness and diversity of the vaginal microbiota and spontaneous preterm birth. *Microbiome* **6**, 117 (2018).
12. Stout, M. J. *et al.* Early pregnancy vaginal microbiome trends and preterm birth. *Am. J. Obstet. Gynecol.* **217**, 356.e1-356.e18 (2017).
13. Hyman, R. W. *et al.* Diversity of the vaginal microbiome correlates with preterm birth. *Reprod. Sci.* **21**, 32–40 (2014).
14. Feehily, C. *et al.* Shotgun sequencing of the vaginal microbiome reveals both a species and functional potential signature of preterm birth. *NPJ Biofilms Microbiomes* **6**, 50 (2020).

15. Kosti, I., Lyalina, S., Pollard, K. S., Butte, A. J. & Sirota, M. Meta-Analysis of Vaginal Microbiome Data Provides New Insights Into Preterm Birth. *Front. Microbiol.* **11**, 476 (2020).
16. Gupta, P., Singh, M. P. & Goyal, K. Diversity of Vaginal Microbiome in Pregnancy: Deciphering the Obscurity. *Front Public Health* **8**, 326 (2020).
17. Ceccarani, C. *et al.* Diversity of vaginal microbiome and metabolome during genital infections. *Sci. Rep.* **9**, 14095 (2019).
18. Chase, A. B., Weihe, C. & Martiny, J. B. H. Adaptive differentiation and rapid evolution of a soil bacterium along a climate gradient. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
19. Zhao, S. *et al.* Adaptive Evolution within Gut Microbiomes of Healthy People. *Cell Host Microbe* **25**, 656-667.e8 (2019).
20. Garud, N. R., Good, B. H., Hallatschek, O. & Pollard, K. S. Evolutionary dynamics of bacteria in the gut microbiome within and across hosts. *PLoS Biol* **17**(1), e3000102 (2019).
21. Garud, N. R. & Pollard, K. S. Population Genetics in the Human Microbiome. *Trends Genet.* **36**, 53–67 (2020).
22. Murovec, B., Deutsch, L. & Stres, B. Computational Framework for High-Quality Production and Large-Scale Evolutionary Analysis of Metagenome Assembled Genomes. *Mol. Biol. Evol.* **37**, 593–598 (2020).
23. Olm, M. R. *et al.* Consistent Metagenome-Derived Metrics Verify and Delineate Bacterial Species Boundaries. *mSystems* **5**(1), e00731-19 (2020).
24. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM: assessing the quality of microbial genomes recovered from isolates, single cells, and metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
25. Vaneechoutte, M. *et al.* Emended description of *Gardnerella vaginalis* and description of *Gardnerella leopoldii* sp. nov., *Gardnerella piovii* sp. nov. and *Gardnerella swidsinskii* sp. nov., with delineation of 13 genomic species within the genus *Gardnerella*. *Int. J. Syst.*

- Evol. Microbiol.* **69**, 679–687 (2019).
26. Hill, J. E., Albert, A. Y. K. & the VOGUE Research Group. Resolution and Cooccurrence Patterns of *Gardnerella leopoldii*, *G. swidsinskii*, *G. piotii*, and *G. vaginalis* within the Vaginal Microbiome. *Infect Immun.* **87**(12), e00532-19 (2019).
  27. Martino, C. *et al.* Context-aware dimensionality reduction deconvolutes gut microbial community dynamics. *Nat. Biotechnol.* **39**, 165–168 (2021).
  28. Mendz, G. L., Petersen, R., Quinlivan, J. A. & Kaakoush, N. O. Potential involvement of *Campylobacter curvus* and *Haemophilus parainfluenzae* in preterm birth. *BMJ Case Rep.* **2014**, (2014).
  29. Suzuki, T. A. & Ley, R. E. The role of the microbiota in human genetic adaptation. *Science* **370**, eaaz6827 (2020).
  30. Ferris, M. J. *et al.* Association of *Atopobium vaginae*, a recently described metronidazole resistant anaerobe, with bacterial vaginosis. *BMC Infect. Dis.* **4**, 5 (2004).
  31. Danielsson, P.-E. Euclidean distance mapping. *Computer Graphics and Image Processing* vol. **14** (3), 227–248 (1980).
  32. Manrique, P., Dills, M. & Young, M. J. The Human Gut Phage Community and Its Implications for Health and Disease. *Viruses* **9**(6), 141 (2017).
  33. Schwebke, J. R., Muzny, C. A. & Josey, W. E. Role of *Gardnerella vaginalis* in the pathogenesis of bacterial vaginosis: a conceptual model. *J. Infect. Dis.* **210**, 338–343 (2014).
  34. Kindschuh, W. F. *et al.* Preterm birth is associated with xenobiotics and predicted by the vaginal metabolome. *Nat Microbiol* 2023 (2023). doi: 10.1038/s41564-022-01293-8.
  35. Weissman, J. L., Hou, S. & Fuhrman, J. A. Estimating maximal microbial growth rates from cultures, metagenomes, and single cells via codon usage patterns. *Proc. Natl. Acad. Sci. U. S. A.* **118**(12), e2016810118 (2021).
  36. Achtman, M. & Wagner, M. Microbial diversity and the genetic nature of microbial species.

- Nat. Rev. Microbiol.* **6**, 431–440 (2008).
37. Bohr, L. L., Mortimer, T. D. & Pepperell, C. S. Lateral Gene Transfer Shapes Diversity of spp. *Front. Cell. Infect. Microbiol.* **10**, 293 (2020).
  38. Hudson, R. R. Linkage disequilibrium and recombination. *Handbook of statistical genetics* (2004).
  39. Kryazhimskiy, S. & Plotkin, J. B. The population genetics of dN/dS. *PLoS Genet.* **4**, e1000304 (2008).
  40. Olm, M. R. *et al.* inStrain profiles population microdiversity from metagenomic data and sensitively detects shared microbial strains. *Nat. Biotechnol.* **39**, 727–736 (2021).
  41. Shenhav, L. & Zeevi, D. Resource conservation manifests in the genetic code. *Science* **370**, 683–687 (2020).
  42. Schloissnig, S. *et al.* Genomic variation landscape of the human gut microbiome. *Nature* **493**, 45–50 (2013).
  43. He, M. *et al.* Evolutionary dynamics of *Clostridium difficile* over short and long time scales. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 7527–7532 (2010).
  44. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for genome-scale analysis of protein functions and evolution. *Nucleic Acids Res.* **28**, 33–36 (2000).
  45. Stokholm, J. *et al.* Antibiotic use during pregnancy alters the commensal vaginal microbiota. *Clin. Microbiol. Infect.* **20**, 629–635 (2014).
  46. Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic resistance database. *Nucleic Acids Res.* **48**, D517–D525 (2020).
  47. Maguire, F. *et al.* Metagenome-assembled genome binning methods with short reads disproportionately fail for plasmids and genomic Islands. *Microb Genom* **6**(10), mgen000436 (2020).
  48. Liao, J. *et al.* Nationwide genomic atlas of soil-dwelling *Listeria* reveals effects of selection

- and population ecology on pangenome evolution. *Nat Microbiol* **6**(8), 1021-1030 (2021).
49. Elovitz, M. A. *et al.* Cervicovaginal microbiota and local immune response modulate the risk of spontaneous preterm delivery. *Nat Commun* **10**, 1305 (2019).
  50. Callahan, B. J. *et al.* Replication and refinement of a vaginal microbial signature of preterm birth in two racially distinct cohorts of US women. *Proc. Natl. Acad. Sci. U. S. A.* **114**, 9966–9971 (2017).
  51. Haque, M. M., Merchant, M., Kumar, P. N., Dutta, A. & Mande, S. S. First-trimester vaginal microbiome diversity: A potential indicator of preterm delivery risk. *Sci. Rep.* **7**, 16145 (2017).
  52. Leung, J. M., Graham, A. L. & Knowles, S. C. L. Parasite-Microbiota Interactions With the Vertebrate Gut: Synthesis Through an Ecological Lens. *Front Microbiol* **9**, 843 (2018).
  53. Morowitz, M. J. *et al.* Strain-resolved community genomic analysis of gut microbial colonization in a premature infant. *Proc Natl Acad Sci U S A* **108**(3), 1128-33 (2011).
  54. Zeevi, D. *et al.* Structural variation in the gut microbiome associates with host health. *Nature* **568**, 43–48 (2019).
  55. Pace, R. M. *et al.* Complex species and strain ecology of the vaginal microbiome from pregnancy to postpartum and association with preterm birth. *Med* (2021) doi:10.1016/j.medj.2021.06.001.
  56. Brown, R. G. *et al.* Vaginal dysbiosis increases risk of preterm fetal membrane rupture, neonatal sepsis and is exacerbated by erythromycin. *BMC Med.* **16**, 9 (2018).
  57. Menard, J. P. *et al.* High vaginal concentrations of *Atopobium vaginae* and *Gardnerella vaginalis* in women undergoing preterm labor. *Obstet. Gynecol.* **115**, 134–140 (2010).
  58. Kumar, S. *et al.* The Vaginal Microbial Signatures of Preterm Birth Delivery in Indian Women. *Front. Cell. Infect. Microbiol.* **11**, 622474 (2021).
  59. Yuan, W., Chen, L. & Bernal, A. L. Is elevated maternal serum alpha-fetoprotein in the second trimester of pregnancy associated with increased preterm birth risk? *Eur J Obstet*

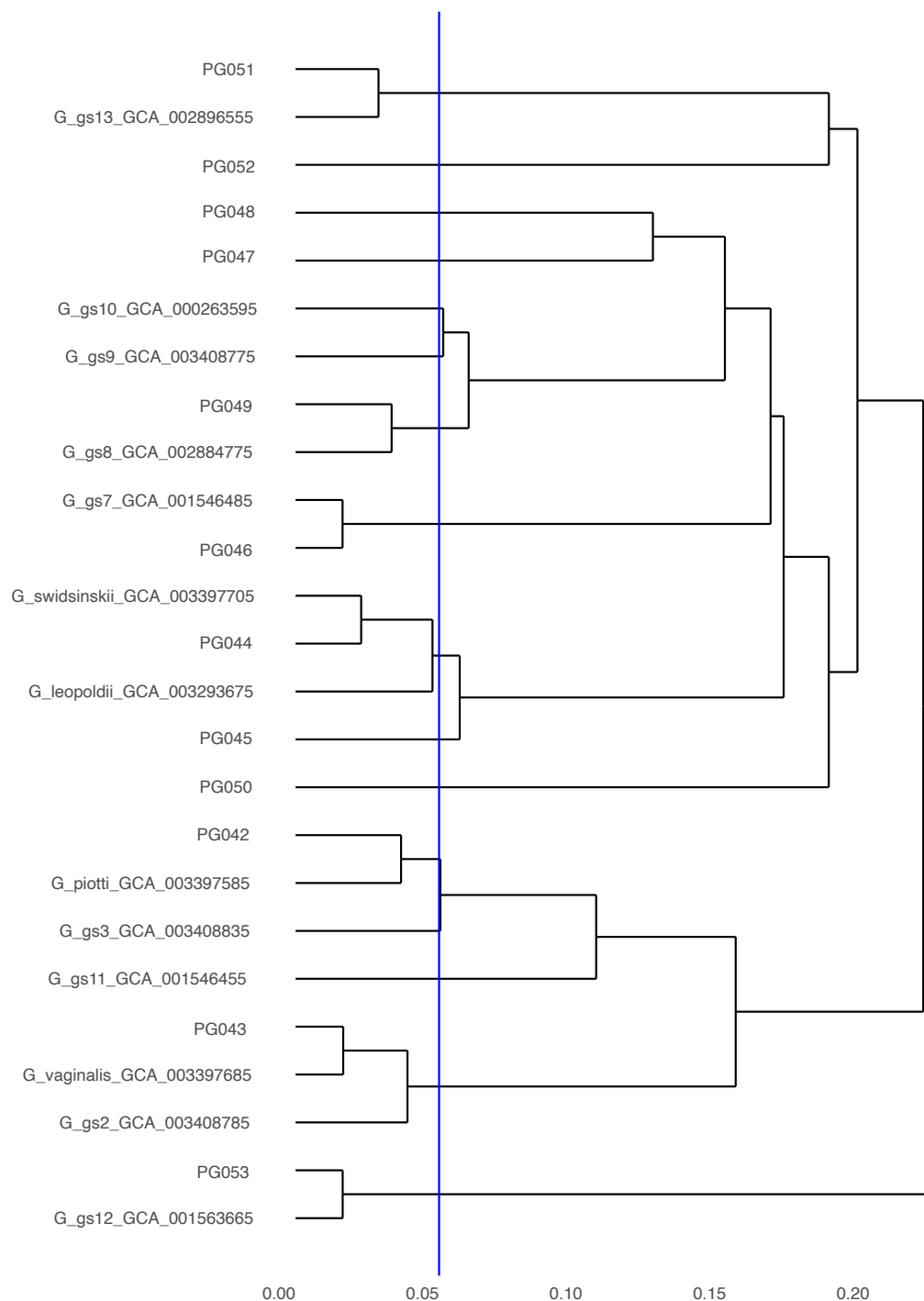
- Gynecol Reprod Biol.* **45**(1), 57-64 (2009).
60. Simhan, H. N., Caritis, S. N., Krohn, M. A. & Hillier, S. L. Elevated vaginal pH and neutrophils are associated strongly with early spontaneous preterm birth. *Am. J. Obstet. Gynecol.* **189**, 1150–1154 (2003).
  61. Otto, S. P. & Barton, N. H. The evolution of recombination: removing the limits to natural selection. *Genetics* **147**, 879–906 (1997).
  62. Cooper, T. F. Recombination speeds adaptation by reducing competition between beneficial mutations in populations of *Escherichia coli*. *PLoS Biol.* **5**, e225 (2007).
  63. Martinez-Gutierrez, C. A. & Aylward, F. O. Strong Purifying Selection Is Associated with Genome Streamlining in Epipelagic Marinimicrobia. *Genome Biol. Evol.* **11**, 2887–2894 (2019).
  64. Gerson, K. D. *et al.* A non-optimal cervicovaginal microbiota in pregnancy is associated with a distinct metabolomic signature among non-Hispanic Black individuals. *Sci. Rep.* **11**, 22794 (2021).
  65. Gharthey, J., Bastek, J. A., Brown, A. G., Anglim, L. & Elovitz, M. A. Women with preterm birth have a distinct cervicovaginal metabolome. *Am. J. Obstet. Gynecol.* **212**, 776.e1-776.e12 (2015).
  66. Terzic, M. *et al.* Periodontal Pathogens and Preterm Birth: Current Knowledge and Further Interventions. *Pathogens* **10**, 730 (2021).
  67. da Costa, A. C. *et al.* Identification of bacteriophages in the vagina of pregnant women: a descriptive study. *BJOG* **128**, 976–982 (2021).
  68. Bargaza, R. A. & Cunha, B. A. Aminoglycosides in gynecology. *Int. Urogynecol. J.* **3**, 197–207 (1992).
  69. Amstey, M. S. Chloramphenicol therapy in pregnancy. *Clinical infectious diseases: an official publication of the Infectious Diseases Society of America* vol. 30 237 (2000).
  70. Rick, A.-M. *et al.* Group B Streptococci Colonization in Pregnant Guatemalan Women:



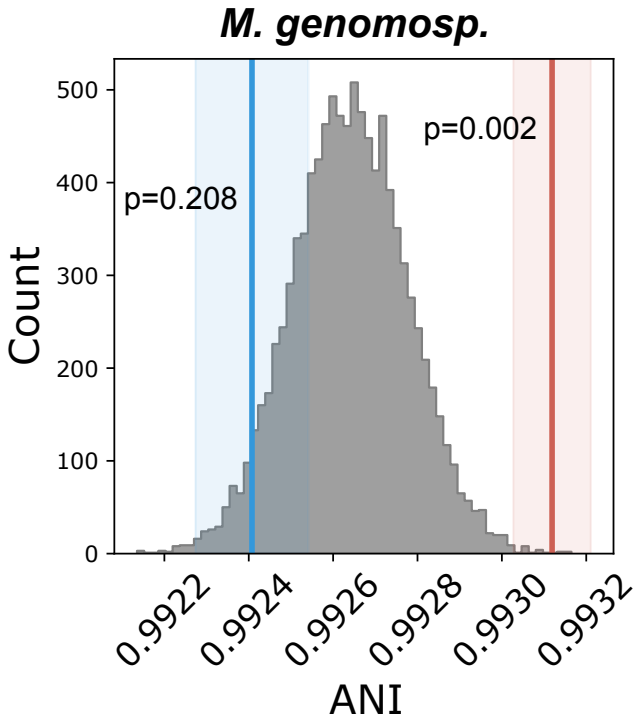
- Prevalence, Risk Factors, and Vaginal Microbiome. *Open Forum Infect Dis* **4**, ofx020 (2017).
71. Rosenbaum, P. R. & Rubin, D. B. The central role of the propensity score in observational studies for causal effects. *Biometrika* **70**, 41–55 (1983).
  72. Gálvez, E. J. C. *et al.* Distinct Polysaccharide Utilization Determines Interspecies Competition between Intestinal *Prevotella* spp. *Cell Host Microbe* **28**, 838-852.e6 (2020).
  73. Yang, S., Liebner, S., Svenning, M. M. & Tveit, A. T. Decoupling of microbial community dynamics and functions in Arctic peat soil exposed to short term warming. *Mol. Ecol.* **30**, 5094–5104 (2021).
  74. Kieser, S., Zdobnov, E. M. & Trajkovski, M. Comprehensive mouse microbiota genome catalog reveals major difference to its human counterpart. *PLoS Comput. Biol.* **18**, e1009947 (2022).
  75. Chevalier, C. *et al.* Warmth Prevents Bone Loss Through the Gut Microbiota. *Cell Metab.* **32**, 575-590.e7 (2020).
  76. Kieser, S., Brown, J., Zdobnov, E. M., Trajkovski, M. & McCue, L. A. ATLAS: a Snakemake workflow for assembly, annotation, and genomic binning of metagenome sequence data. *BMC Bioinformatics* **21**, 257 (2020).
  77. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
  78. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
  79. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
  80. Kang, D. *et al.* MetaBAT 2: an adaptive binning algorithm for robust and efficient genome reconstruction from metagenome assemblies. *PeerJ* **7**, e7359 (2019).
  81. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate

- genomic comparisons that enables improved genome recovery from metagenomes through de-replication. *ISME J.* **11**, 2864–2868 (2017).
82. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site identification. *BMC Bioinformatics* **11**, 119 (2010).
  83. Huerta-Cepas, J. *et al.* eggNOG 5.0: a hierarchical, functionally and phylogenetically annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic Acids Res.* **47**, D309–D314 (2019).
  84. Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**(2), giab008 (2021).
  85. Chaumeil, P.-A., Mussig, A. J., Hugenholtz, P. & Parks, D. H. GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database. *Bioinformatics* (2019) doi:10.1093/bioinformatics/btz848.
  86. Darriba, D., Taboada, G. L., Doallo, R. & Posada, D. ProtTest 3: fast selection of best-fit models of protein evolution. *Bioinformatics* **27**, 1164–1165 (2011).
  87. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
  88. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree display and annotation. *Nucleic Acids Res.* **49**, W293–W296 (2021).
  89. Nies, L. de *et al.* PathoFact: A pipeline for the prediction of virulence factors and antimicrobial resistance genes in metagenomic data. *Microbiome* **9**(1), 49 (2021).
  90. Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: A practical and powerful approach to multiple testing. *J. R. Stat. Soc.* **57**, 289–300 (1995).

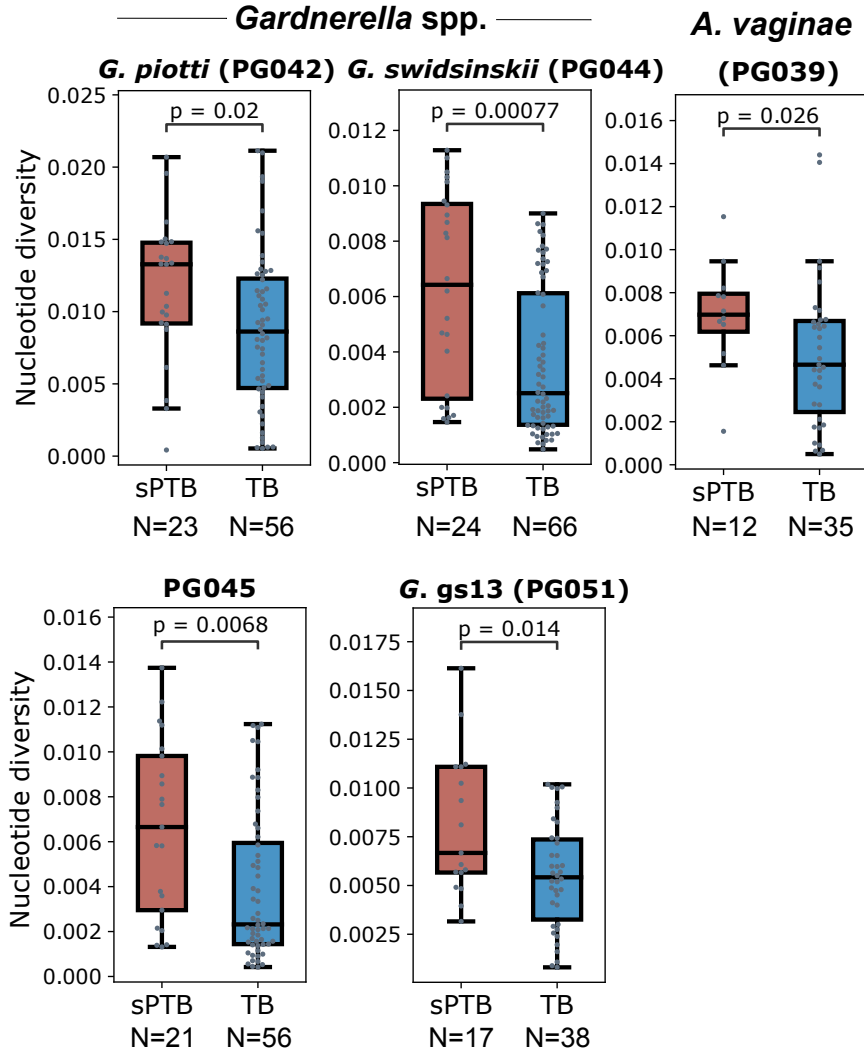
## Supplementary figures



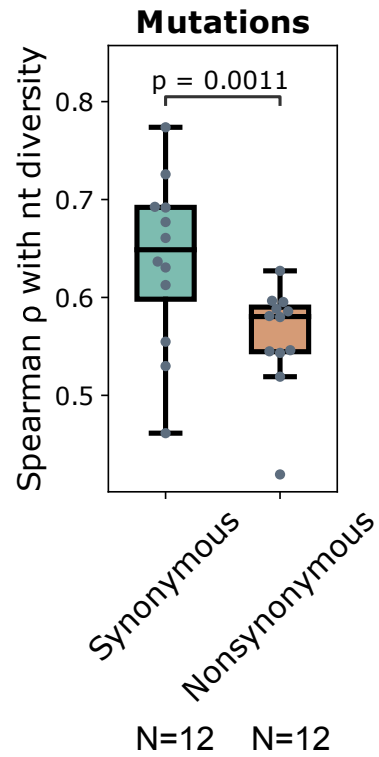
**Supplementary Fig. 1 | Dendrogram of representative MAGs annotated as *G. vaginalis* (PG42-53) and reference genomes of 13 *Gardnerella* species defined in ref. <sup>25</sup> based on ANI. Blue line shows an ANI value of 0.95, which is a cutoff for prokaryotes species.**



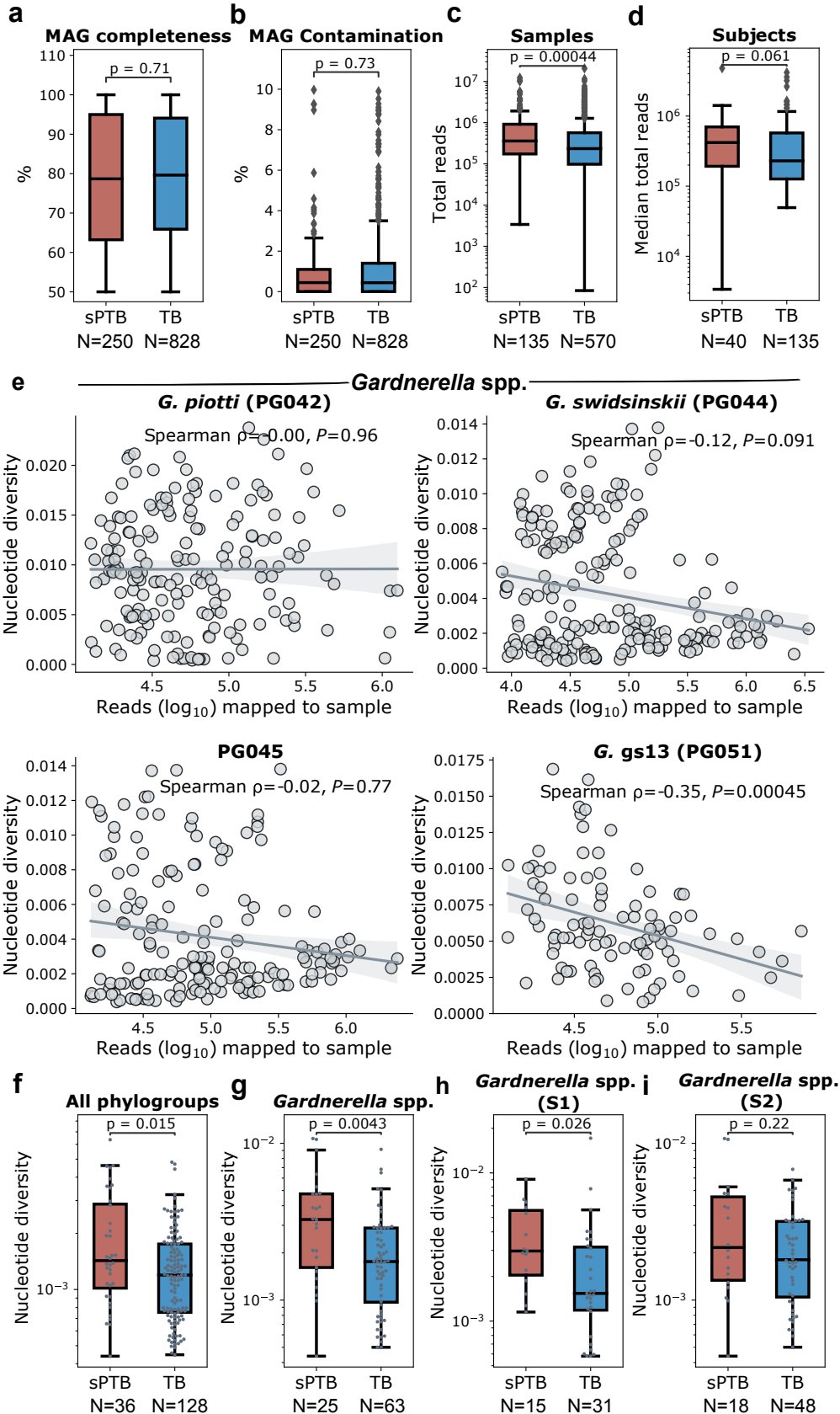
**Supplementary Fig. 2 | The distribution of average inter-host nucleotide identity (ANI) of MAGs classified as *M. genomosp.*** The gray histogram illustrates the null distribution. The red and blue line and shaded area indicate the average value and standard deviation of ANI observed in sPTB and in TB, respectively, calculated from 10,000 repetitions;  $p$ , significance.



**Supplementary Fig. 3 | Nucleotide diversity of vaginal microbial populations.** Median genome-wide nucleotide diversity along pregnancies of four *Gardnerella* spp. phylogroups and a phylogroup classified as *A. vaginae*, compared between sPTB and TB. Box, IQR; line, median; whiskers, 1.5\*IQR;  $p$ , two-sided Mann-Whitney U.



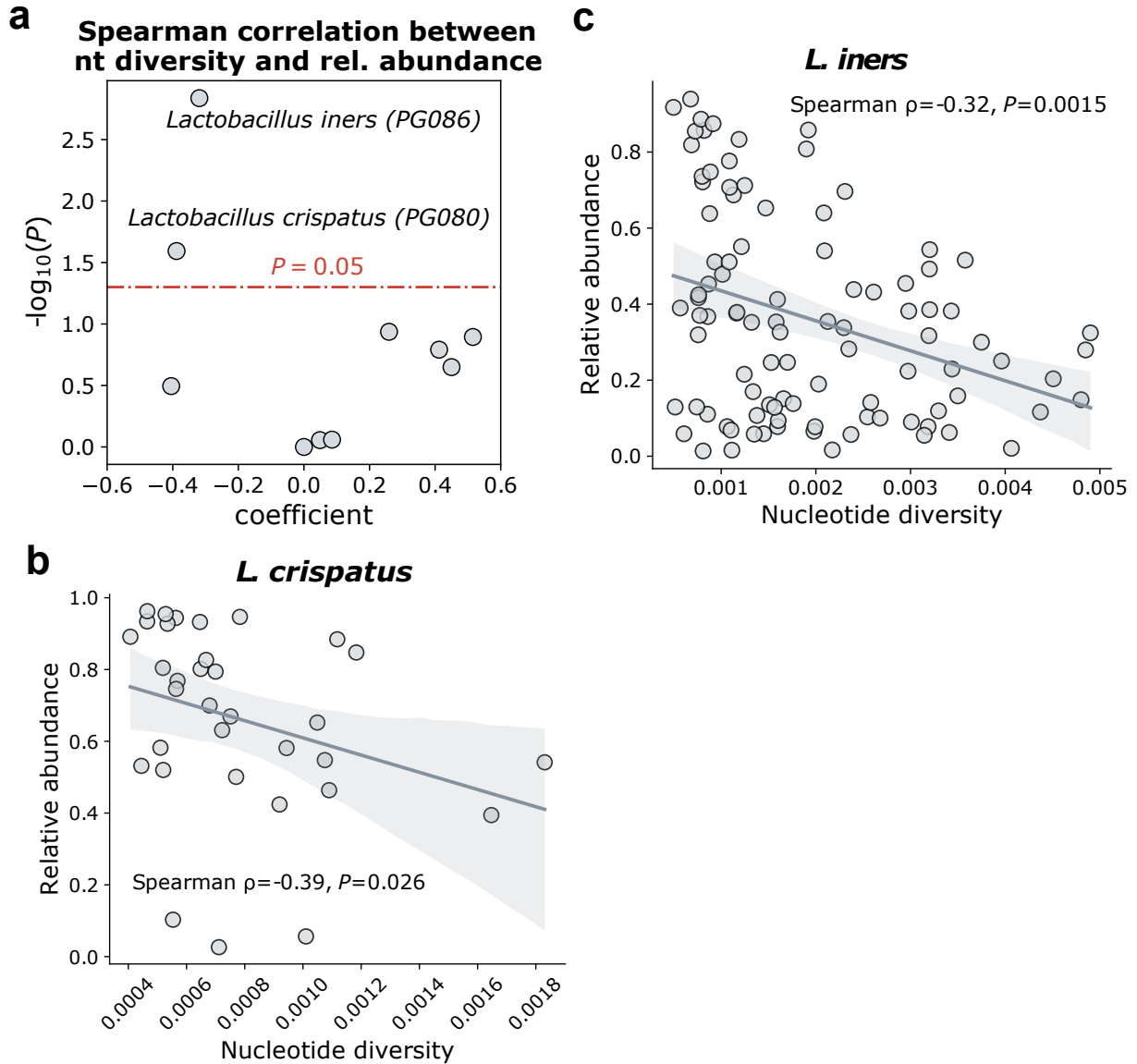
**Supplementary Fig. 4 | Spearman correlation coefficient between nucleotide diversity and number of synonymous mutations and nonsynonymous mutations of genes across 12 *Gardnerella* spp. phylogroups.** Median spearman correlation coefficient of genes of each phylogroup based on median gene nucleotide diversity along pregnancy, compared between synonymous mutations and nonsynonymous mutations. Box, IQR; line, median; whiskers, 1.5\*IQR;  $p$ , two-sided Student T test.



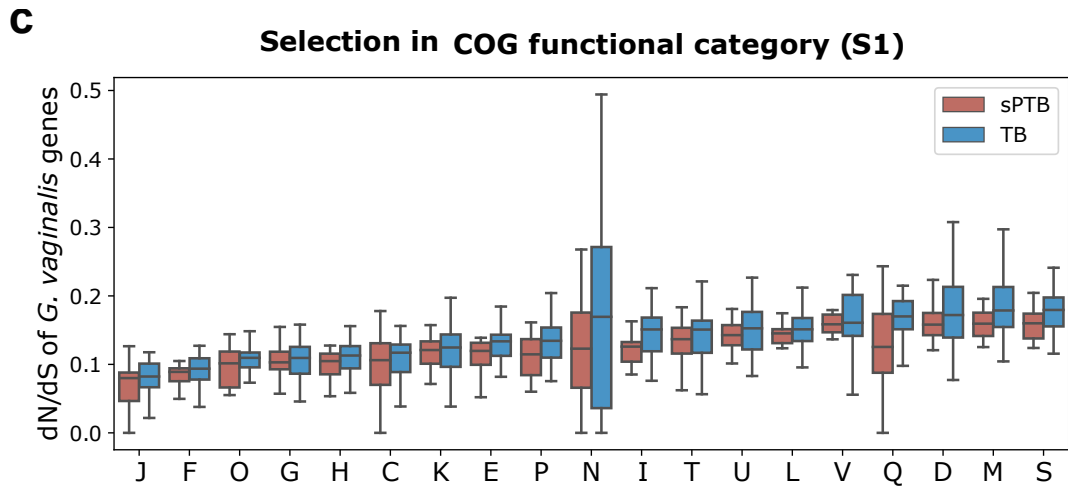
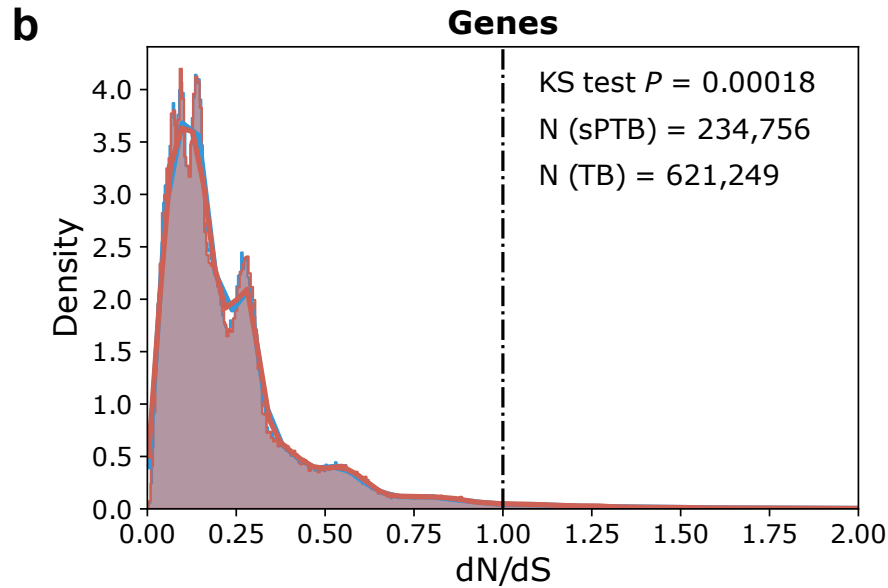
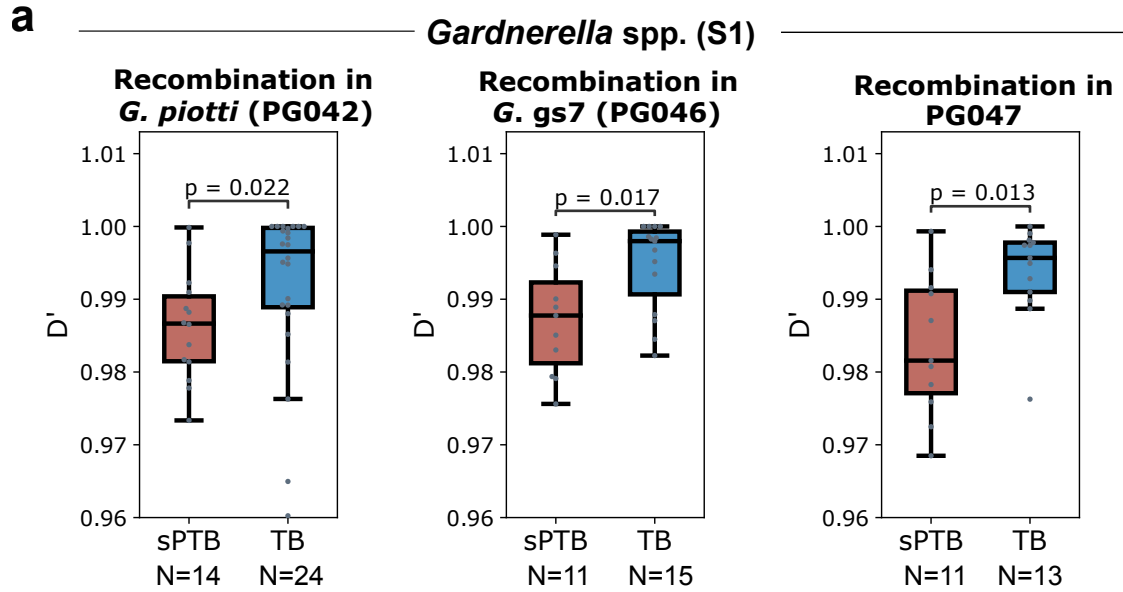
**Supplementary Fig. 5 | The association between microdiversity and sPTB is not biased by**

**sequencing depth. a,b**, Completeness (a) and contamination (b) of MAGs, compared between sPTB and TB. **c,d**, Total read counts of samples (c) and median read count (d) along each pregnancy compared between sPTB and TB. **e**, Spearman correlation between genome-wide nucleotide diversity and reads mapped to each of the four *Gardnerella* spp. phylogroups that show difference in nucleotide diversity between sPTB and TB in **Fig. S2a**. The line and the shaded area depict the best-fit trendline and the 95% confidence interval (mean  $\pm$  1.96 s.e.m.) of the linear regression. **f,g**, Median genome-wide nucleotide diversity along pregnancy of all phylogroups (f) and *Gardnerella* spp. (g), compared between sPTB and TB based on  $10^5$  reads sampled from each sample. **h,i**, Median genome-wide nucleotide diversity of *Gardnerella* spp. along the first (S1, h) and second (S2, i) halves of pregnancy, compared between sPTB and TB based on  $10^5$  reads sampled from each sample. Box, IQR; line, median; whiskers,  $1.5 \times$  IQR; *p*, two-sided Mann-Whitney.

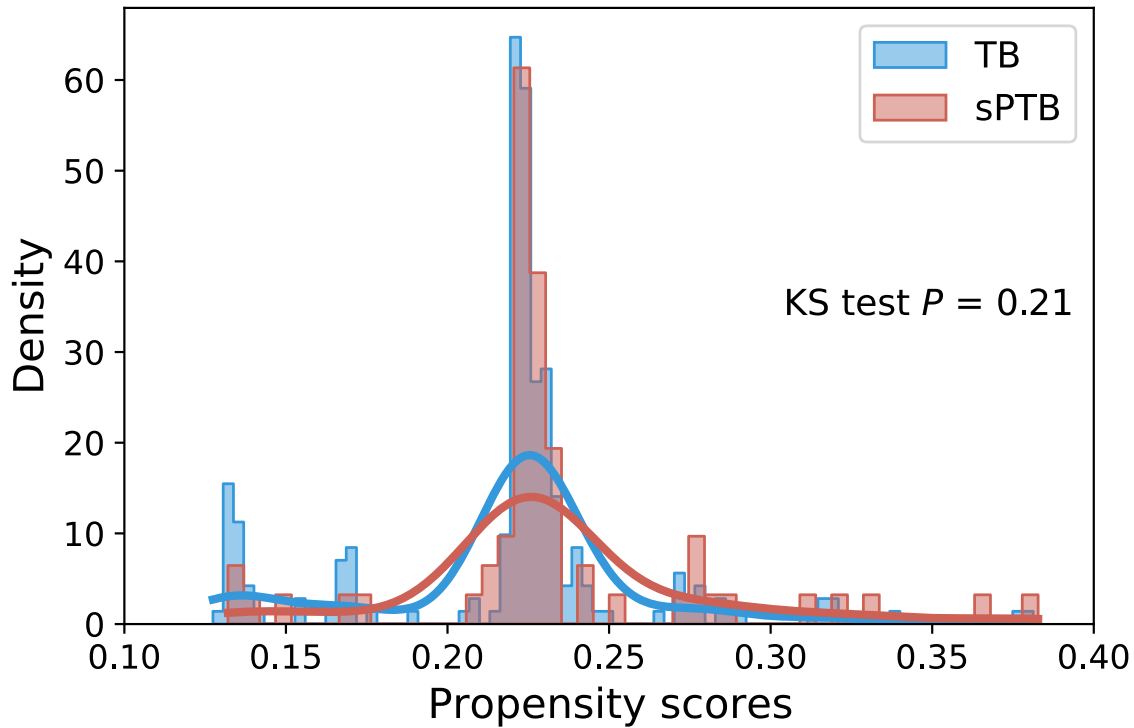




**Supplementary Fig. 6 | Spearman correlation between genome-wide nucleotide diversity and relative abundance of non-*Gardnerella* phylogroups.** a. Volcano plot illustrating the Spearman correlation (significance, y-axis; coefficient, x-axis) between median genome-wide nucleotide diversity and relative abundance along pregnancies. Phylogroups above the red dashed line have a  $P < 0.05$ . b,c, Spearman correlation between median genome-wide nucleotide diversity and relative abundance of *L. crispatus* (b) and *L. iners* (c) along pregnancy. The line and the shaded area depict the best-fit trendline and the 95% confidence interval (mean  $\pm$  1.96 s.e.m.) of the linear regression.



**Supplementary Fig. 7 Evolutionary forces on the vaginal microbiome. a.** Median  $D'$  along the first half of pregnancy (S1) of *Gardnerella* spp. phylogroups compared between sPTB and TB. Lower  $D'$  indicates more frequent recombination. Box, IQR; line, median; whiskers,  $1.5 \times \text{IQR}$ ;  $P$ , two-sided Mann-Whitney. **b.** Density of median (along pregnancy) of dN/dS of genes in sPTB (red) and in TB (blue).  $P$ : Kolmogorov–Smirnov (KS) test. **c.** dN/dS of *Gardnerella* spp. genes compared between sPTB and TB by COG functional categories, displayed for the first (S1, **c**) half of pregnancy. C, Energy production and conversion; D, Cell cycle control, cell division, chromosome partitioning; E, Amino acid transport and metabolism; F, Nucleotide transport and metabolism; G, Carbohydrate transport and metabolism; H, Coenzyme transport and metabolism; I, Lipid transport and metabolism; J, Translation, ribosomal structure and biogenesis; K, Transcription; L, Replication, recombination and repair; M, Cell wall/membrane/envelope biogenesis; N, Cell Motility, O: Post-translational modification, protein turnover, chaperones; P, Inorganic ion transport and metabolism; Q, Secondary metabolites biosynthesis, transport and catabolism; T, Signal transduction mechanisms; U, Intracellular trafficking, secretion, and vesicular transport; V, Defense mechanisms.



**Supplementary Fig. 8 | Distribution of propensity scores of women groups based on income, age, and race using a logistic regression model.** The histogram is smoothed using a kernel. sPTB: red, TB: blue,  $P$ : Kolmogorov–Smirnov (KS) test.

## **Supplementary tables**

**Supplementary Table 1** Genome assembly features of representative MAGs for phylogroups and taxonomy.

**Supplementary Table 2** eggNOG functional annotation of genes.