1  **Complementing 16S rRNA gene amplicon sequencing with estimates of total bacterial load to**

2  **infer absolute species concentrations in the vaginal microbiome**

3

4  Florencia Tettamanti Boshier[1], Sujatha Srinivasan[1], Anthony Lopez[1], Noah G. Hoffman[2], Sean

5  Proll[5], David N. Fredricks [1,3,4,5], Joshua T. Schiffer[1,3,4]

6

7  [1]Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, USA

8  [2]Department of Laboratory Medicine, University of Washington, Seattle, USA

9  [3]Clinical Research Division, Fred Hutchinson Cancer Research Center, Seattle, USA

10  [4]Department of Medicine, University of Washington, Seattle, USA

11  [5]Department of Microbiology, University of Washington, Seattle, USA

12  [6]Department of Biostatistics,  University of Washington, Seattle, USA

12    Whereas 16S rRNA gene amplicon sequencing quantifies relative abundances of bacterial taxa,

13    variation in total bacterial load between samples restricts its ability to reflect absolute

14    concentration of individual species. Quantitative PCR (qPCR) can quantify individual species, but

15    it is not practical to develop a suite of qPCR assays for every bacterium present in a diverse

16    sample. We analyzed 1320 samples from 20 women with a history of frequent bacterial vaginosis,

17    who self-collected vaginal swabs daily over 60 days. We inferred bacterial concentrations by

18    taking the product of species relative abundance (assessed by 16S rRNA gene amplicon

19    sequencing) and total bacterial load (measured by broad-range 16S rRNA gene qPCR). $Log_{10}$-

20    converted inferred concentrations correlated with targeted qPCR (r = 0. 935, p<2.2e-16) for seven

21    key bacterial species. The mean inferred concentration error varied across bacteria, with rarer

22    bacterial vaginosis-associated bacteria associated with larger errors. 92% of errors >0.5 $log_{10}$

23    occurred when relative abundance was <10%. Many errors occurred during early bacterial

24    expansion or late contraction. When relative abundance of a species is >10%, inferred

25    concentrations are reliable proxies for targeted qPCR. However, targeted qPCR is required to

26    capture bacteria at low relative abundance, particularly with BV-associated bacteria during the

27    early onset of bacterial vaginosis.

28

29 **Introduction**

30 For most infectious diseases, the absolute concentration of a single pathogen is often the most

31 specific marker of disease severity and therapeutic response(1–3). In constrast, studies of

32 bacterial communities usually rely on broad-range consensus sequence PCR of taxonomically

33 informative genes (such as 16S rRNA) coupled with next generation sequencing (NGS) to assess

34 relative, but not absolute abundances of bacteria. At a mechanistic level, specific combinations

35 of bacteria and bacterial gene products are thought to play a causative role in the pathogenesis

36 of many microbiome associated conditions(4–6), and this approach of characterizing the

37 microbiota is valuable. However, absolute concentration of individual bacterial taxa within

38 communities may be a better predictor of biological activity or disease risk compared to relative

39 abundances of these taxa. Quantitating absolute concentration of individual species with qPCR

40 is time intensive, requires generation of a standard curve for each organism using known

41 concentrations of DNA, is expensive and only available in specialized laboratories. Moreover,

42 each qPCR assay requires significant development and validation costs. qPCR is therefore not

43 typically comprehensive for all species in a community. Moreover, selection of the most

44 appropriate species for analysis may reflect investigator bias.

45 A method to infer absolute concentration of multiple bacterial species from NGS data would be

46 extremely useful for the field including studies of the vaginal microbiome. NGS amplicon

47 sequencing is a fractional approach that has been used to help define conditions such as bacterial

48 vaginosis (7–10), and to identify enhanced risk for other sexually transmitted infections and pre-

49 term delivery (11,12). However, total bacterial load may vary significantly between and within

50 individuals over time even over the course of a single day (8). Therefore, relative abundances

51   may not accurately represent absolute concentrations. Consequently, as shown recently in the

52   gut microbiome, relative abundances may identify spurious disease associations which may in

53   fact be driven by total microbial load (13).

54   Here, we demonstrate that multiplying relative abundance data (composition) by estimates of

55   total bacterial DNA as measured by a broad-range 16S rRNA gene qPCR assay provides useful

56   estimates of absolute concentrations of bacterial DNA for a given species. These inferred

57   concentrations have already been used in studies of the penile microbiome, though without

58   formal validation (14). Herein we validate inferred concentrations by comparison of absolute

59   concentrations measured by targeted qPCR assay for seven key species in the vaginal

60   microbiome.  We find that whereas inferred concentrations are accurate for most samples, they

61   are prone to error when relative abundance is low and may misrepresent kinetics of individual

62   species during critical periods of expansion and clearance.

63

64   **Matherials and Methods**

65   *Ethics statement.* Vaginal samples were collected using protocol 417, which was approved by the

66   institutional review board (IRB) at the University of Washington (approval no.: STUDY00000398).

67   All participants provided written informed consent prior to study enrollment. Consent forms

68   were approved by the IRB as part of protocol 417.

69   *Study Population.* The study population was comprised of 20 women enrolled in a longitudinal

70   study of bacterial vaginosis (BV) natural history at the University of Washington Virology

4

71    Research clinic between 2015 and 2017. At enrollment, participants were given sufficient swabs

72    for three times daily swabs over 60 days. Diagnosis, sample collection, storage, and processing

73    of swabs are as described in (15). Participants were also given a study diary to record symptoms

74    of BV, antibiotic use, menstruation, sexual activity and other medical events. In total, as some

75    participants occasionally skipped samples, we analyzed 1320 data points for each of the seven

76    key species.

77    ***DNA Extraction and Quantitative Polymerase Chain Reaction (qPCR).*** DNA was extracted from

78    vaginal swabs using the BiOstic Bacteremia DNA Isoaltion Kit (Mobio, Carlsbad, CA). Sham swab

79    without human contact were extracted in parallel to assess contamination from reaction buffers

80    or the collection swabs. No template water controls were included to determine if there was any

81    contamination from PCR reagents. Each sample was evaluated for PCR inhibition (Khot et. al.

82    BMC Infectious Diseases. 2008) and total bacterial concentrations in each sample were measured

83    using a qPCR assay that targets the V3-V4 region of the 16S rRNA gene (Srinivasan et al. PloS ONE

84    2012). Concentrations of specific vaginal bacteria were measured using qPCR assays targeting 7

85    key vaginal bacteria: *Atopobium vaginae*, BV-associated bacterium 2 (BVAB2), *Gardnerella*

86    *vaginalis, Lactobacillus crispatus, Lactobacillus jensenii, Lactobacillus iners,* and *Megasphaera*

87    (combined species 1 and 2) species (12,16,17). We measured relative abundances of bacterial

88    taxa using broad-range PCR targeting the V3-V4 region of the 16S rRNA gene with next-

89    generation sequencing on the Illumina MiSeq instrument (Illumina, San Diego, CA)(18). The

90    *DADA2* pipeline was used to infer sequene variants from raw reads for subsequent analysis (19).

91    Sequences were classified using the phylogenetic placement tool *pplacer* (20) and a curated

92    reference set of vaginal bacteria (8). In subsequent text, we use NGS to to refer to data generated

93    using broad-range PCR and sequencing. Sequence reads have been submitted to the NCBI Short

94    Read Archive (in submission. Accession numbers pending). Relative abundances and absolute

95    concentrations of specific vaginal bacteria were measured on all samples in two participants and

96    in daily morning samples for the remaining 18. We performed qPCR on all samples collected from

97    each participant, but for the purpose of this work only consider the morning samples.

98    All data generated or analysed during this study are included in the supplementary material.

99    ***Statistical considerations.*** We calculated inferred concentrations using equation 1.

100   *Equation 1*

101   $$\text{Inferred Concentration} = \text{Relative abundance} \times \text{Total Bacterial Load}$$
102   $$\text{(16S rRNA gene copies)} \quad \text{(\%)} \quad \text{(16S rRNA gene copies)}$$

103   We present the plots and related calculations on a $\log_{10}$ scale. To keep all values finite, zero

104   relative abundance (%) were mapped to 1e-5 and zero inferred concentrations were mapped to

105   1. The choice of this mapping changes some of the numerical results presented here, namely the

106   correlation coefficient and the clustering class of the samples.    However, the general

107   observations being made are consistent.

108   We defined the error of inferred concentration, IC error, as

109   *Equation 2*

110   $$\text{IC error} = \log_{10}(\text{absolute concentration}) - \log_{10}(\text{inferred concentration})$$
111
112   Rates of change per day where calculated between any two consecutive time points which were

113   18-36 hours apart. Rates were calculated from $\log_{10}$ converted values for relative abundance and

114   inferred and absolute concentration. We defined the error in rates from inferred concentrations,

6

115 rIC error, as

116
117 *Equation 3*

118    $\text{rIC error} = \text{rates(absolute concentration)} - \text{rates(inferred concentration)}$

119

120 Comparison of means was done using the t.test function in R (21). We used Pearson's correlation

121 coefficient for all correlation analysis. This was done using the cor.test function in the stats

122 package in R (21). Correlation coefficients were compared using the Cocor package in R (22). The

123 suite provides 10 test for overlapping correlations, i.e. measurements taken from the same data

124 set. All test were significant, but we report the value of the Hittner test here for simplicity.

125 The Breusch-Pagan test was used to test the heteroskedasticity of the linear regression model of

126 the relative abundance and inferred concentration vs absolute concentration. It tests whether

127 the variance of the erros from a regression is dependent on the values of the independent

128 variables. This was implemented using the bptest of the lmtest package in R (23).

129 We constructed the dendrograms for clustering analysis by complete linkage hierarchical

130 clustering of species abundance and/or concentration based on Euclidean distance between all

131 sample pairs.

132 We tested concordance between pairs of dendrograms using the entanglement coefficient found

133 in the dendextend package in R (24). To calculate the coefficient, first all the samples are

134 numbered in the order they appear for each tree. The coefficient is then calculated by taking the

135 Euclidean distance of these two vectors which is then normalised by the worst case entanglement

136 value (i.e. the Euclidean distance when the order of the two dendrograms is opposite). The

137 entanglement coefficient thus defined ranges from 0 to 1, with 0 indicating perfect alignment

138    between the dendorgrams and 1 a complete mismatch.

139

140 **Results**

141 ***Bacterial kinetics in 20 women with frequent recurrent BV.*** The bacterial kinetics observed for a

142 single participant are shown in **Figure 1a and b.** The individual shown underwent dynamic

143 changes in bacterial profile with notable shifts between low to high diversity states. The bacterial

144 kinetics of 19 other participants can be found **in Figure S1**. As previously noted, high diversity

145 states were often concurrent with high absolute concentrations of *Gardnerella vaginalis*,

146 *Atopobium vaginae*, BVAB2 and *Megasphaera,* which all have been associated with bacterial

147 vaginosis (BV) (8,10,25).

148 In 5 of the participants, shifts in composition appear less abruptly when measured by

149 single species qPCR than by NGS. For example, for the participant shown in Figure 1, the absolute

150 concentration of *A. vaginae* increases on day 17 (hour 415), but its relative abundance does not

151 show a consistent increase until day 28 (hour 671) although there are some non-zero abundances

152 in 4/9 samples before this point. From day 0 to day 28 (hour 168), the participant received

153 metronidazole for BV: qPCR shows an exponential decline in BV-associated species absolute

154 concentrations(26); yet, NGS shows a much more abrupt shift towards *L. iners* predominance.

155 NGS can also fail to capture low-level colonization of bacteria, such as that of *G. vaginalis* on days

156 12 to 11 ( hours 150 and 261). Several high-diversity samples have highly prevalent species which

157 were not measured with qPCR in this study, such as *Prevotella bivia* from day 28 onwards (hour

158 671). These observations, which can be made for many of the individuals in this cohort, highlight

159 that qPCR provides more granular estimates for measuring single species kinetics while NGS is

160 optimal to estimate bacterial diversity in high diversity communities.

9

161

***Relative abundances may misclassify single species absolute concentration due to shifts in total***

***bacterial load.*** We compared absolute concentration and relative abundance from the same

samples measured within individuals over the course of the study. Examples for two species , *L.*

*crispatus* and *Megasphaera*, are shown in **Figure 2a** and **b** (examples for the remaining five

species are in supplementary **Figure S2**). There were time points in which absolute and relative

abundance measures demonstrated opposing or differing kinetics, often due to concurrent large

shifts in total bacterial load. These are indicated by arrows in **Figure 2a** and **b**. Thus, relative

abundance may misrepresent absolute concentration when not accounting for total bacterial

load.

171

***Inferred concentrations are predictive of absolute concentrations measured by qPCR.*** For each

species we calculated inferred concentrations by multiplying total bacterial load by NGS-relative

abundance as shown in equation 1. We then compared these with absolute concentration as

measured by targeted qPCR assay for the seven key species. For each species, inferred bacterial

concentration closely tracked absolute concentration for most samples **(dotted line in Figure 2a**

and **b and Figure S2)**. In many instances and for most species there were no obvious extreme

discordance noted (**Figure 2a** and **S2**). For some species however, such as *Megasphaera* and

BVAB2, inferred concentration consistently overestimated absolute concentration by an order of

magnitude (**Figure 2b** and **S2d**). In a subset of samples, for all species, inferred concentration was

zero while qPCR levels were positive leading to profound discordance between inferred and

absolute concentration: this was most often noted at low absolute concentration **(Figure 2a and**

10

183   **b)**.

184       We compared correlation between relative abundance and absolute concentration

185   (r=0.936, P < 2.2e-16 **Figure 3a**) to correlation between inferred concentration and absolute

186   concentration (r=0.935, p<2.2e-16 **Figure 3b**). The two correlation coefficients are not

187   statistically different (Hittner test, p>0.08)(22).  Species-specific correlations were noted. For

188   inferred concentrations, *Megasphaera* and BVAB2 produced the strongest correlation followed

189   by *L. crispatus, A. vaginae* and *L. jensenii; G. vaginalis* and *L. iners*, which are often present at

190   moderate concentrations (~$10^6$ 16S rRNA gene copies), had the weakest correlations though

191   correlations coefficients for all species were high **(Table 1)**.

192       We defined error of inferred concentration, IC error, as in equation 2.  While there was a

193   large range in errors for non-zero inferred concentrations (**Figure 3b,** range: -7.32 log10 (16S

194   rRNA gene copies) – 2.66 log10(16S rRNA gene copies)), the mean IC error (-0.319 log10(16S rRNA

195   gene copies)) and standard deviation (0.999 log10(16S rRNA gene copies)) were low. Moreover,

196   the median IC error for most species approximated zero with samples within the interquartile

197   range demonstrating minimal IC error (**Figure 4a**). However, for BVAB2 and *Megasphaera*, the

198   interquartile range of IC error, while narrow, was all less than zero, implying consistent

199   overestimation of absolute concentration by inferred concentration (pair-wise t-test p<0.05).

200   There was a trend towards global underestimation of *G. vaginalis* using inferred concentration

201   **(Figure 4a)**.

202

11

203    ***Low relative abundance is the major source of IC error.*** The variance in the relationship with

204    absolute concentration tended to be inversely proportional to species concentrations [Breusch-

205    Pagan test; p = 0.06] highlighting that a larger range of IC errors tended to be reported at lower

206    bacterial loads (**Figure 3a**). Accordingly, 93% of >0.5 IC errors were accounted for by relative

207    abundances below 10% and 85% by relative abundances below 1%. Many of these IC errors

208    occurred on double negatives – samples for which inferred concentration was zero and absolute

209    concentration was reported at threshold. When these samples were removed from the analysis,

210    84% of >0.5 IC errors were accounted for by relative abundances <10% and 66% by relative

211    abundances below 1% **(Figure 4b).** The median absolute concentration above the limit of

212    detection for >0.5 IC errors was 5.95 log10(16S rRNA gene copies) (IQR: 4.03 – 7.88 , range: 1.97

213    – 10.39 ).

214          We defined false positive samples as non-zero inferred concentration values when

215    absolute concentration qPCR values were at or below the detection threshold and false negatives

216    as zero-values for inferred concentration when absolute concentrations were above the

217    detection threshold. False negatives were more common (23.6% of samples) than false positives

218    (3.17 % of samples) which demonstrates that targeted qPCR is more sensitive for single species

219    detection than NGS.

220          The incidence of false negatives was not equal across species, with *G. vaginalis* having the

221    highest percentage of false negatives, followed by *L. inners* and *A. vaginae* (*L.crispatus* 13.8%, *L.*

222    *jensenii* 31.1 %, *L. iners* 35.1%, *G. vaginalis* 60.4%, *A. vaginae* 35.3%, *Megasphaera* 5.40%, BVAB2

223    9.84%). The higher percentages of false negative for some species occurred because they are

224 often present at moderate concentrations, near the relative abundance error threshold. The

225 median qPCR value for false negative samples was 3.92 $\log_{10}$ (16S rRNA gene copies) (IQR: 2.88 –

226 4.82, range: 1.97 -7.84 ), again showing that IC errors generally occur at lower bacterial loads.

227    Total bacterial load measured by broad-range qPCR assay was frequently below the sum

228 of the concentration of all seven species measured by targeted qPCR assays (37.6% per species

229 per sample). Non-zero inferred concentrations from samples with underestimates of total

230 bacterial load consistently overpredicted absolute concentration (one-tailed t-test p<2.6e-4) and

231 did so more than at other points (pair-wised t-test p<2.2e-16) (**Figure 4c**). Non-zero inferred

232 concentrations from samples with known underestimates of total bacterial load had a median IC

233 error of 0.171 log10 (16S rRNA gene copies) (IQR  -0.138 - 0.447, range -7.31 – 2.66) compared

234 to -0.368  $\log_{10}$ (16S rRNA gene copies)  (IQR -0.638 - -0.143, range: -6.54 – 1.42) in other samples.

235    *L. crispatus* had the highest percentage of false positives (*L.crispatus* 8.42%, *L. jensenii*

236 1.08%, *L. iners* 3.56%, *G. vaginalis* 0.46%, *A. vaginae* 3.07%, *Megasphaera* 1.12%, BVAB2 1.79%).

237 The median relative abundance of false positives across all samples was extremely low 0.06%

238 (IQR: 0.04 – 0.11%, range 0.0007 – 36.8%).

239 ***Concentrations inferred from NGS predicts observed absolute concentration regardless of***

240 ***sample diversity or sequencing depth.*** Inferred concentrations did not disproportionally record

241 misleading results from low or high diversity samples as measured by Shannon Diversity index

242 (**Figure 5a**). Moreover, we observed occasional large absolute IC errors across all sequencing

243 depths (**Figure 5b**). Low bacterial abundance was the primary source of absolute IC error

244 regardless of diversity or sequencing depth (**Figure 5a**  and **b)**. Larger than 0.5 absolute IC error

13

245  was observed across all raw species counts, but the largest absolute IC errors (above >2) were

246  almost exclusively associated with raw species counts below 100 **(Figure 5c)**.

247

248  ***Inferred concentration estimates are predictive of most temporal changes in single species***

249  ***bacterial load.*** We examined whether inferred concentration is a useful tool for evaluating

250  individual species kinetics by determining changes in bacterial levels over the course of a day.

251  Rates of change in relative abundances correlated only weakly with absolute concentrations

252  [r=0.271, p<2.2e-16]. Moreover, 17.1% of the time, we observed a change in relative abundance

253  in the opposite direction to that of absolute concentration (top-left and bottom-right quadrants

254  in **Figure 6a**). This type of error occurred commonly for both the most abundant (e.g. *L. crispatus*)

255  and rarer species (e.g. BVAB2).

256        Rates of change in inferred concentration showed improved correlation with rates of

257  change in absolute concentration [pmcc=0.392, p<2.2e-16]. The mean rIC error (defined in the

258  **Methods**) was low (-2.71 x10$^{-3}$, SD: 1.54 log$_{10}$(16S rRNA gene copies) per hour) though the range

259  of rIC errors was high (-9.29 − 9.31 log$_{10}$(16S rRNA gene copies) per hour), indicating occasional

260  samples with very poor prediction. Inferred concentrations decreased the sign rIC error rate by

261  more than 50% (from 17.1% to 7.97%, **Figure 6b**).

262        **Figure 7a** shows a typical profile of *A. vaginae* absolute levels and sample-to-sample

263  change, to demonstrate the two types of rIC errors which were most common to the data. The

264  first were large positive or negative rates which occurred when one of two consecutive points

265  had an inferred concentration of zero, while absolute concentration was detectable by qPCR.

266  These points resulted in dramatic overestimation of growth or contraction rates for individual

14

267 species across all samples (**Figure 6b & 7b,** right upper and left lower quadrants). Such rIC error

268 often occurred when species were transitioning to or from a low abundance ($<10^6$ 16S rRNA gene

269 copies per sample). The second type of rIC error, occurred when two consecutive points had

270 inferred concentrations of zero, resulting in underestimation of growth or contraction rates for

271 individual species **(Figure 7b)**. This phenomenon also commonly occurred when a species was

272 transitioning to or from a low abundance ($<10^6$ 16S rRNA gene copies per sample). These two

273 forms of transitions accounted for 91.7% of rIC error > 0.05 **(Figure 7c)**. If all transitions involving

274 a zero value were eliminated from the analysis, we observed excellent correlation between

275 inferred and observed rate of change (r=0.876, p<2.2e-16). It follows that inferred concentrations

276 do not capture kinetics during microbial blooming or contraction when bacteria are at low

277 concentration or not detected using the less sensitive broad-range PCR with NGS approach.

278 However, inferred concentrations can be used to estimate individual species growth and

279 contraction rates when bacteria are present at higher concentrations such as $>10^6$ 16S rRNA gene

280 copies/swab .

281

282 ***Complete linkage clustering by inferred and absolute concentrations shows general agreement.***

283 To assess whether inferred concentrations provide similar or disparate classification of samples,

284 we clustered samples using complete linkage hierarchical clustering based on Euclidean distances

285 (21) by inferred and absolute concentrations **(Figure S3)**. We compared the resulting

286 dendrograms using the entanglement coefficient from the dendextend package in R (24), where

287 a value of 1 corresponds to complete discordance and a value of 0 indicates perfect alignment.

15

288    The two dendrograms were found to be in agreement, with a low entanglement coefficient 0.11.

289

290        We next determined the number of clusters using NbClust package in R (27). Absolute

291    concentration identifies two whereas inferred concentration identifies three clusters. The third

292    cluster arose from a general disctinction between samples dominated by *L. crispatus* from *L. iners*

293    as the inferred concentrations had a lower threshold (1 copy per swab) than the qPCR (93.8

294    copies per swab).

295

296    ***Inferred concentration may provide the most comprehensive overview of individual species***

297    ***kinetics.*** Inferred concentrations can be calculated for all species captured by NGS. In **Figure 1**

298    and **S1**, we show the inferred concentrations of the most abundant species across all samples.

299    We imposed a 1% relative abundance threshold to limit the possible 0.5 IC error described in

300    **Figure 4b**. This relative abundance cut-off results in abrupt appearance and disapparence of

301    organisms. Although we cannot validate our projections for species outside of the seven key

302    bacterial species for which we have targeted qPCR assays, inferred concentrations have the

303    potential to describe the kinetics of relevant species present at moderate to high concentrations

304    during bacterial shifts in the microbiome.

305        We carried out complete linkage hierarchical clustering based on Euclidean distance by

306    inferred concentration and relative abundance for the 20 most abundant species of the data set

307    **(Figure S4)**. The resulting dendrograms showed general agreement, with an entanglement

308    coefficient of 0.12. Both techniques identified two clusters defined by high concentration *G.*

309    *vaginalis* and high diversity versus *Lactobacillus* predominance (27).

310 **Discussion**

311 An ideal assay that characterizes bacterial communities in an ecological niche would capture

312 several metrics including species composition, diversity, and quantity as reflected by absolute

313 concentration of all species present. Broad-range PCR of phylogenetically informative genes

314 followed by NGS is the most commonly used approach and captures the first two metrics.

315 However, because total bacterial levels may shift dramatically over narrow time intervals, relative

316 abundance measures by NGS do not reflect absolute concentration. While it is possible to

317 circumvent this issue with targeted (taxon specific) qPCR, these assays are expensive, time

318 consuming and only available in specialized laboratories. Invariably, the absolute concentration

319 of many relevant species is left unmeasured due to these constraints.

320 This measurement gap is of high relevance to clinical studies of the human microbiome in which

321 total bacterial load may not be stable. It is biologically plausible that the absolute levels of critical

322 species are more predictive of health and disease states than relative levels, as is the case with

323 classical single pathogen infectious diseases. Moreover, serial measurements of absolute levels

324 are necessary to fully capture non-linear microbial dynamic changes which relate to inter-species

325 competition for limited resources.

326 Using a large longitudinal dataset of the vaginal microbiome notable for frequent changes

327 between low and high diversity states, we demonstrate that the absolute concentration of a

328 given species can be inferred by multiplying the total bacterial quantity by its relative abundance

329 as measured by NGS. Given that quantitating total bacterial load is affordable and available to

330 many laboratories, this simple approach may allow estimation of absolute concentration without

17

331    needing to perform qPCR on all samples.

332    Our technique is remarkably predictive of absolute concentration with certain key exceptions.

333    Species such as BVAB2 and *Megasphaera* which were often present at low absolute

334    concentration were notable for high precision but slight inaccuracy: inferred concentration

335    consistently slightly overestimated the abundance for these species. This result highlights that

336    individual comparisons between inferred and absolute concentration must be performed for all

337    species of interest. Other than in an exploratory fashion, we do not advocate the use of inferred

338    concentration for species which have not been validated in depth with targeted qPCR assays and

339    compared to absolute concentration.

340    Second, our approach has a very high IC error rate when relative abundance is low, or zero. In

341    our qPCR dataset low level colonization of certain species often precedes a surge in levels prior

342    to this species predominating. Because qPCR is more sensitive than NGS for low amounts of

343    bacterial DNA, and because inferred concentration relies on NGS, inferred concentration will

344    often miss persistent low-level colonization, as well as the critical early growth phase or late

345    contraction phase of relevant species. Despite this fact, inferred concentration performs

346    remarkably well at estimating growth and decay rates at the single species level, provided these

347    rates are estimated based on positive sequential samples. One might be able to improve accuracy

348    of the inferred concentrations by increasing sequencing depth or improving the accuracy of

349    measurements of the total bacterial load.

350    A final issue not addressed by our technique is the limitation inherent to comparing bacterial

351    quantities between species using qPCR based on differing amplification efficiencies of different

352   assays. This variability may arise from different bacterial targets having varying GC content,

353   secondary structures and amplification product size.  In this sense, absolute concentration by

354   qPCR may not be a perfect gold standard for comparing inferred concentration.

355   In summary, we developed and validated a simple, user-friendly method to estimate absolute

356   species abundance in complex polymicrobial communities. This method is best employed when

357   species are present at >10% relative abundance and must be validated for each species of

358   interest. Ultimately, inferred concentration of one or several species may serve as a more

359   predictive variable of disease association, compared to relative abundance, and may advance our

360   understanding of how specific environmental and host factors influence microbial

361   concentrations.

362   **Acknowledgements**

365   **Author Contributions**

366   J.T.S, S.S. and D.N.F conceived and designed the experiments.  A.L. performed the experiments.

367   N.G.H. managed the NGS bioinformatic pipeline.  S.P. managed data integration and contributed

368   to figure generation.  J.T.S. and F.A.T.B. conceived the idea of inferred concentration.  F.A.T.B

369   completed the analysis, contributed to figure generation and wrote the manuscript.

370   **Competing Interests**

371   The authors declare no competing interests.

19

372 **References**

373 1. Falade-Nwulia OO, Naggie S, Nahass RG, Kim AY, Scott JD, Ghany MG, et al. Hepatitis C

374 Guidance 2018 Update: AASLD-IDSA Recommendations for Testing, Managing, and

375 Treating Hepatitis C Virus Infection. Clin Infect Dis. 2018;67(10):1477–92.

376 2. File TM. Highlights from international clinical practice guidelines for the treatment of

377 acute uncomplicated cystitis and pyelonephritis in women: A 2010 update by the

378 infectious diseases society of America and the european society for microbiology and

379 infectious . Infect Dis Clin Pract. 2011;19(4):282–3.

380 3. Saag MS, Benson CA, Gandhi RT, Hoy JF, Landovitz RJ, Mugavero MJ, et al. Antiretroviral

381 drugs for treatment and prevention of HIV infection in adults: 2018 recommendations of

382 the international antiviral society-USA panel. JAMA - J Am Med Assoc. 2018;320(4):379–

383 96.

384 4. Bisgaard H, Hermansen MN, Buchvald F, Loland L, Halkjaer LB, Bonnelykke K, et al. of the

385 Airway in Neonates. N Engl J Med. 2007;357:1487–95.

386 5. Dejea CM, Fathi P, Craig JM, Boleij A, Taddese R, Geis AL, et al. Patients with familial

387 adenomatous polyposis harbor colonic biofilms containing tumorigenic bacteria. Science

388 (80- ). 2018;359(6375):592–7.

389 6. Costello SP, Hughes PA, Waters O, Bryant R V., Vincent AD, Blatchford P, et al. Effect of

390 Fecal Microbiota Transplantation on 8-Week Remission in Patients with Ulcerative Colitis:

391 A Randomized Clinical Trial. JAMA - J Am Med Assoc. 2019;321(2):156–64.

392 7. Srinivasan S, Morgan MT, Fiedler TL, Djukovic D, Hoffman NG, Raftery D, et al. Metabolic

393 signatures of bacterial vaginosis. MBio. 2015;6(2):1–16.

394   8.    Srinivasan S, Hoffman NG, Morgan MT, Matsen FA, Fiedler TL, Hall RW, et al. Bacterial

395         communities in women with bacterial vaginosis: High resolution phylogenetic analyses

396         reveal relationships of microbiota to clinical criteria. PLoS One. 2012;7(6).

397   9.    Ravel J, Gajer P, Abdo Z, Schneider GM, Koenig SSK, McCulle SL, et al. Vaginal microbiome

398         of reproductive-age women. Proc Natl Acad Sci [Internet].

399         2011;108(Supplement_1):4680–7. Available from:

400         http://www.pnas.org/cgi/doi/10.1073/pnas.1002611107

401   10.   Gajer P, Brotman RM, Bai G, Sakamoto J, Schütte UME, Zhong X, et al. Temporal

402         dynamics of the human vaginal microbiota. Sci Transl Med [Internet].

403         2012;4(132):132ra52. Available from:

404         http://stm.sciencemag.org/content/scitransmed/4/132/132ra52.full

405   11.   Nelson DB, Hanlon A, Nachamkin I, Haggerty C, Mastrogiannis DS, Liu C, et al. Early

406         pregnancy changes in bacterial vaginosis-associated bacteria and preterm delivery.

407         Paediatr Perinat Epidemiol. 2014;28(2):88–96.

408   12.   McClelland RS, Lingappa JR, Srinivasan S, Kinuthia J, John-Stewart GC, Jaoko W, et al.

409         Evaluation of the association between the concentrations of key vaginal bacteria and the

410         increased risk of HIV acquisition in African women from five cohorts: a nested case-

411         control study. Lancet Infect Dis [Internet]. 2018;18(5):554–64. Available from:

412         http://dx.doi.org/10.1016/S1473-3099(18)30058-6

413   13.   Vandeputte D, Kathagen G, D'Hoe K, Vieira-Silva S, Valles-Colomer M, Sabino J, et al.

414         Quantitative microbiome profiling links gut community variation to microbial load.

415         Nature [Internet]. 2017;551(7681):507–11. Available from:

416         http://dx.doi.org/10.1038/nature24460

417   14.    Liu CM, Prodger JL, Tobian AAR, Abraham AG, Price LB. crossm Penile Anaerobic

418         Dysbiosis as a Risk Factor for HIV Infection. :1–10.

419   15.    Mayer BT, Matrajt L, Casper C, Krantz EM, Corey L, Wald A, et al. Dynamics of persistent

420         oral cytomegalovirus shedding during primary infection in ugandan infants. J Infect Dis.

421         2016;214(11):1735–43.

422   16.    Fredricks DN, Fiedler TL, Thomas KK, Mitchell CM, Marrazzo JM. Changes in vaginal

423         bacterial concentrations with intravaginal metronidazole therapy for bacterial vaginosis

424         as assessed by quantitative PCR. J Clin Microbiol. 2009;47(3):721–6.

425   17.    Srinivasan S, Liu C, Mitchell CM, Fiedler TL, Thomas KK, Agnew KJ, et al. Temporal

426         variability of human vaginal bacteria and relationship with bacterial vaginosis. PLoS One.

427         2010;5(4).

428   18.    Garcia K, Celustka K, Srinivasan S, Loeffelholz T, Fiedler TL, Aker S, et al. Stool Microbiota

429         at Neutrophil Recovery Is Predictive for Severe Acute Graft vs Host Disease After

430         Hematopoietic Cell Transplantation. Clin Infect Dis. 2017;65(12):1984–91.

431   19.    Callahan BJ, McMurdie PJ, Rosen MJ, Han AW, Johnson AJA, Holmes SP. DADA2: High-

432         resolution sample inference from Illumina amplicon data. Nat Methods. 2016;13(7):581–

433         3.

434   20.    FA M, RB K, EV A. pplacer: linear time maximum-likelihood and Bayesian phylogenetic

435         placement of sequences onto a fixed reference tree. BMC Bioinformatics [Internet].

436         2010;11:538. Available from: http://dx.doi.org/10.1186/1471-2105-11-538

437   21.    Team RC. R: A Language and Environment for Statistical Computing. Vienna, Austria: R

438        Foundation for Statistical Computing; 2018.

439    22.    Diedenhofen B, Musch J. Cocor: A comprehensive solution for the statistical comparison

440        of correlations. PLoS One [Internet]. 2015;10(4):1–12. Available from:

441        http://dx.doi.org/10.1371/journal.pone.0121945

442    23.    Achim Zeileis TH. Diagnostic Checking in Regression Relationships. R News [Internet].

443        2010;2(3):7–10. Available from: http://cran.r-project.org/doc/Rnews/

444    24.    Sieger T, Hurley CB, Fiser K, Beleites C. Interactive Dendrograms: The *R* Packages **idendro**

445        and **idendr0**. J Stat Softw [Internet]. 2017;76(10). Available from:

446        http://www.jstatsoft.org/v76/i10/

447    25.    Fredricks DN, Fiedler TL, Marrazzo JM. Molecular Identification of Bacteria Associated

448        with Bacterial Vaginosis. N Engl J Med. 2005;353(18):1899–911.

449    26.    Marrazzo JM, Fiedler TL, Srinivasan S, Mayer BT, Schiffer JT, Fredricks DN. Rapid and

450        Profound Shifts in the Vaginal Microbiota Following Antibiotic Treatment for Bacterial

451        Vaginosis. J Infect Dis. 2015;212(5):793–802.

452    27.    Charrad M, Ghazzali N, Boiteau V, Niknafs A. Package 'NbClust.' 2015;9. Available from:

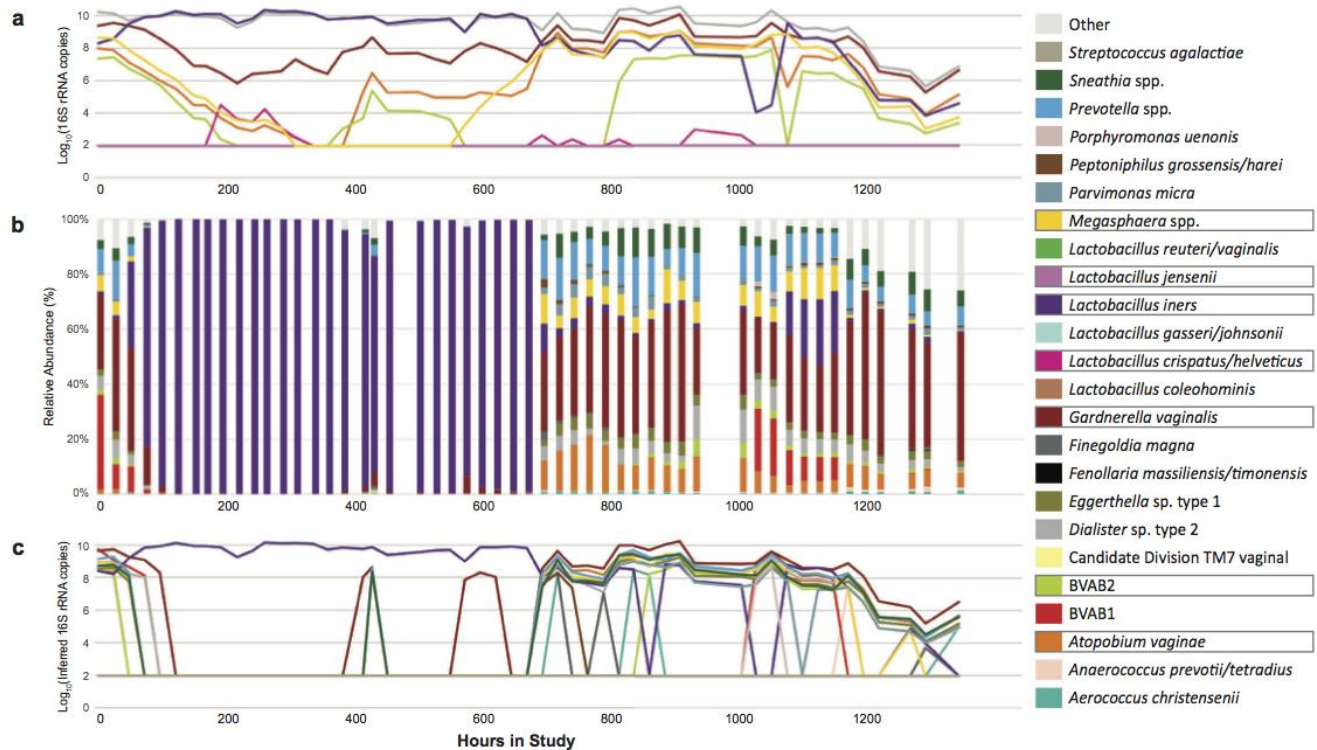453        https://cran.r-project.org/web/packages/NbClust/NbClust.pdf

454

**Figure 1. Complex bacterial kinetics in the vaginal niche in one representative study participant.** Daily samples from a woman, Participant 18, who performed self-swabbing of the vagina were analyzed by: **a)** targeted qPCR of seven specific species, **b)** high throughput sequencing using 16S rRNA and **c)** inferred concentration for species with relative abundance above 1%. qPCR allows measures of absolute concentration, whereas broad range PCR with sequencing provides a measure of bacterial diversity in a given sample. Targeted qPCR often detects shifts in single species prior to NGS. Inferred concentration follows qPCR more closely than relative abundance does and may project concentration of species for which targeted qPCR assays are not available.

**Figure 2. Relative abundances estimates can misrepresent actual concentrations due to shifts in total bacterial load.** Examples of species-specific profiles in two participants for two different species **a)** *L. crispatus*, Participant 06 and **b)** *Megasphaera,* Participant 17. Vertical bars show relative abundance (%, left-y-axis), solid lines are absolute concentrations measured by qPCR and grey line is total bacterial load, dashed lines are inferred concentrations (all right y-axis). The dashed black line indicates detection threshold for qPCR data (93.8 16S rRNA copies). Arrows indicate timepoints when relative abundance changes are discordant from absolute concentration changes, which often occur when bacterial loads shift dramatically, or relative abundance is low. Examples for the remaining species can be found in **Fig S2.**
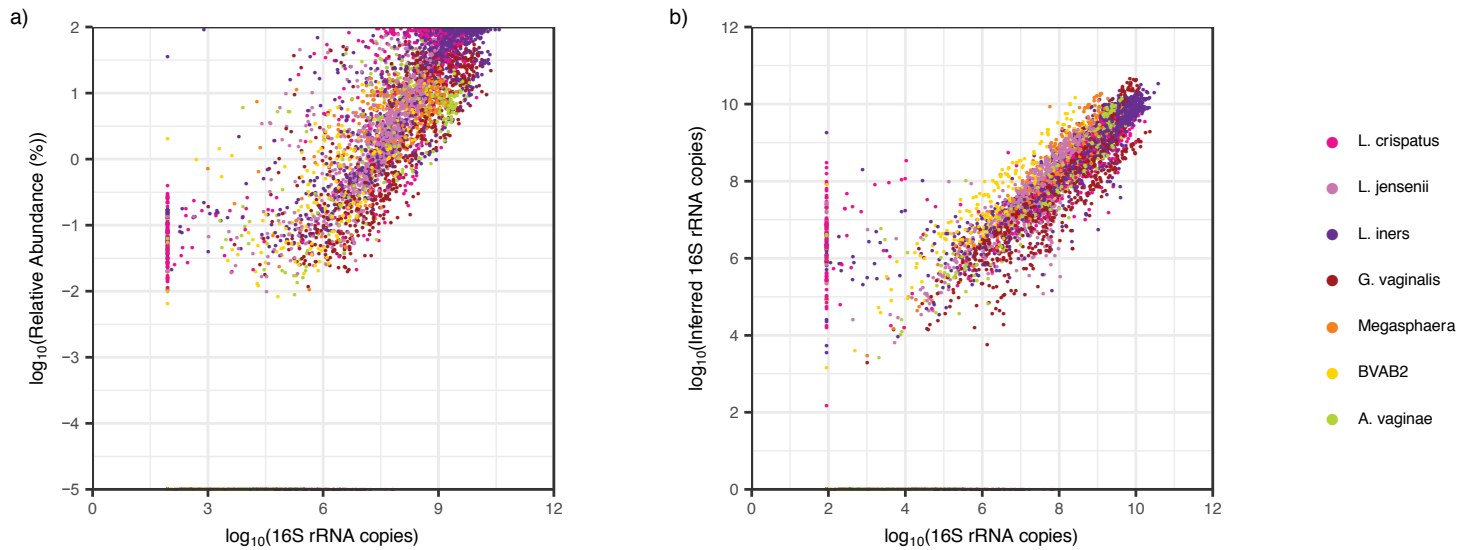
**Figure 3. Inferred concentration correlates more strongly with absolute concentration than relative abundance. a)** Scatter plot of relative abundance vs absolute concentration. Pearson correlation coefficient (pcc), r is 0.936, P<2.2e-16. **b)** Scatter plot of inferred concentration vs absolute concentration. Both axes are plotted on a logarithmic scale. Pearson correlation coefficient (pcc), r is 0.935, P<2.2e-16. Samples which were negative by NGS but not by targeted qPCR are plotted on the x-axis while samples negative by targeted qPCR but positive by NGS are listed on the reported threshold for targeted qPCR ($\log_{10}(93.8)$ 16S rRNA gene copies). Relative abundances and inferred concentrations were generally falsely negative at low absolute concentrations. Variance in the relationship between absolute concentration and relative abundance is inversely proportional to species concentrations (Breusch-Pagan test, P=2e-3). Whereas this relationship was not statistically significant between absolute concentration and inferred abundance (Breusch-Pagan test, P=0.06).

| Species | PCC Relative Abundance | PCC Inferred Abundance |
|---|---|---|
| Megasphaera | 0.969 | 0.978 |
| BVAB2 | 0.942 | 0.952 |
| Lactobacillus crispatus | 0.941 | 0.920 |
| Atopobium vaginae | 0.911 | 0.916 |
| Lactobacillus jensenii | 0.908 | 0.911 |
| Gardnerella vaginalis | 0.881 | 0.890 |
| Lactobacillus iners | 0.863 | 0.872 |

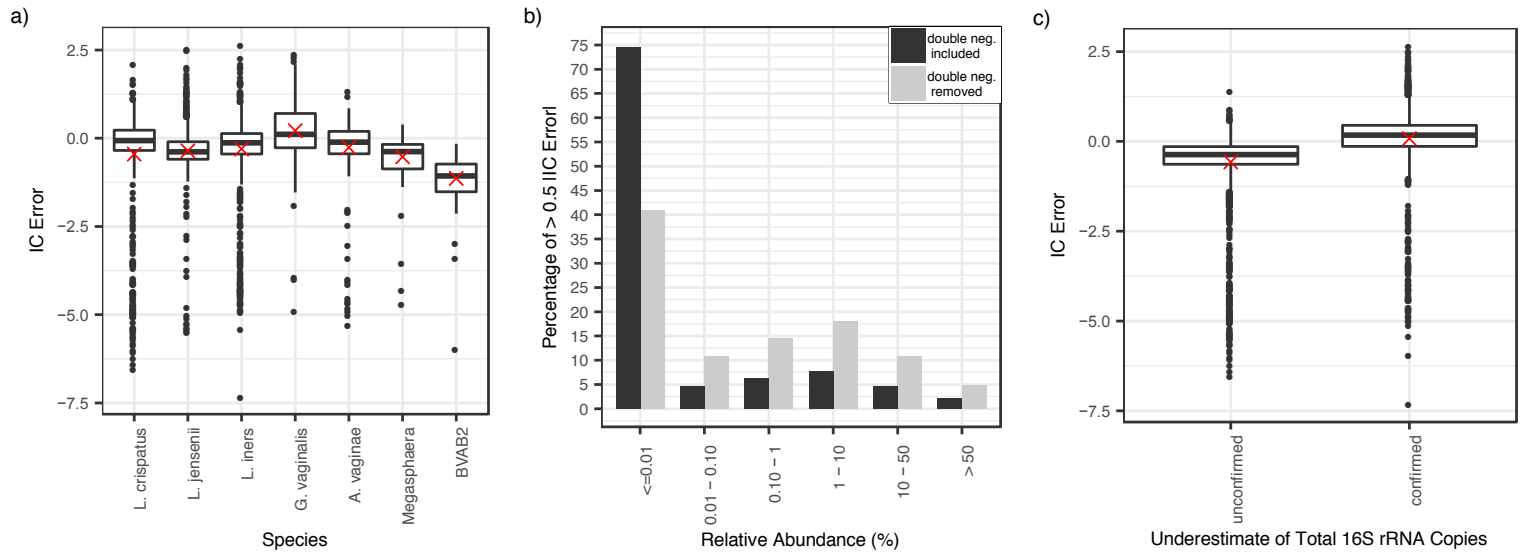**Table 1.** *Pearson correlation coefficients of single species between absolute concentration versus relative abundance and inferred concentration.*

**Figure 4. Low relative bacterial abundance is the major predictor of IC error for inferred concentrations compared to absolute concentrations. a)** Boxplots displaying IC error (equation 2), with zero inferred concentrations removed, indicate low IC error rates overall. Inferred values are consistent overestimates for BVAB2 and *Megasphaera* spp. Boxes are the interquartile range; whiskers are 1.5x the IQR and dots are samples outside of this range; red crosses are mean. **b)** Bar-chart of incidence of >0.5 IC error by relative abundance group. Black is inclusive of double-negatives (0 inferred concentration and threshold absolute concentration): 93% of >0.5 IC errors are accounted for by relative abundances <10% (85% by relative abundances <1%). In grey, concurrent negative samples are removed: 84% of >0.5 IC errors are accounted for by relative abundances <10% (66% by relative abundances <1%). **c)** Boxplots displaying IC error for samples with unconfirmed and confirmed underestimates of total bacterial load by broad range qPCR assay (samples where BR16S is lower than the sum of concentrations of the seven targeted species). Data points with zero inferred concentration were removed. Samples with underestimates of total bacterial load overestimate concentration more than other samples. Overall however, the range of IC error is comparable between both groups. Boxes are the interquartile range; whiskers are 1.5x the IQR and dots are samples outside of this range; crosses are mean.
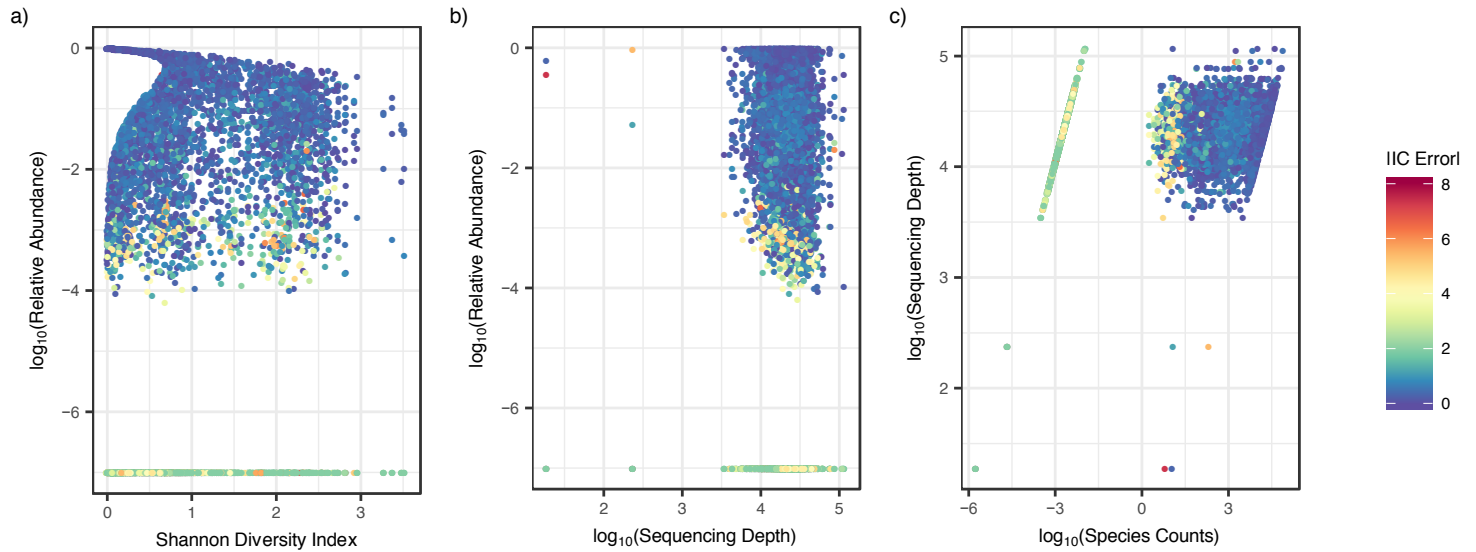
**Figure 5. Sample diversity, sequencing depth and species counts do not impact IC error of inferred concentration.** Scatter-plots color-coded by IC error. Each dot is a sample for a specific species from a single participant. **a)** Relative abundance versus Shannon diversity index. High IC error predominately occurred at low relative abundance but across both low and high diversity samples. **b)** Relative abundance versus sequencing depth. High IC error predominately occurred at low relative abundance but across various levels of sequencing depth. c) Sequencing depth vs species counts. High IC error occurred at species counts below 100, although >0.5 IC error is observed across all species counts.
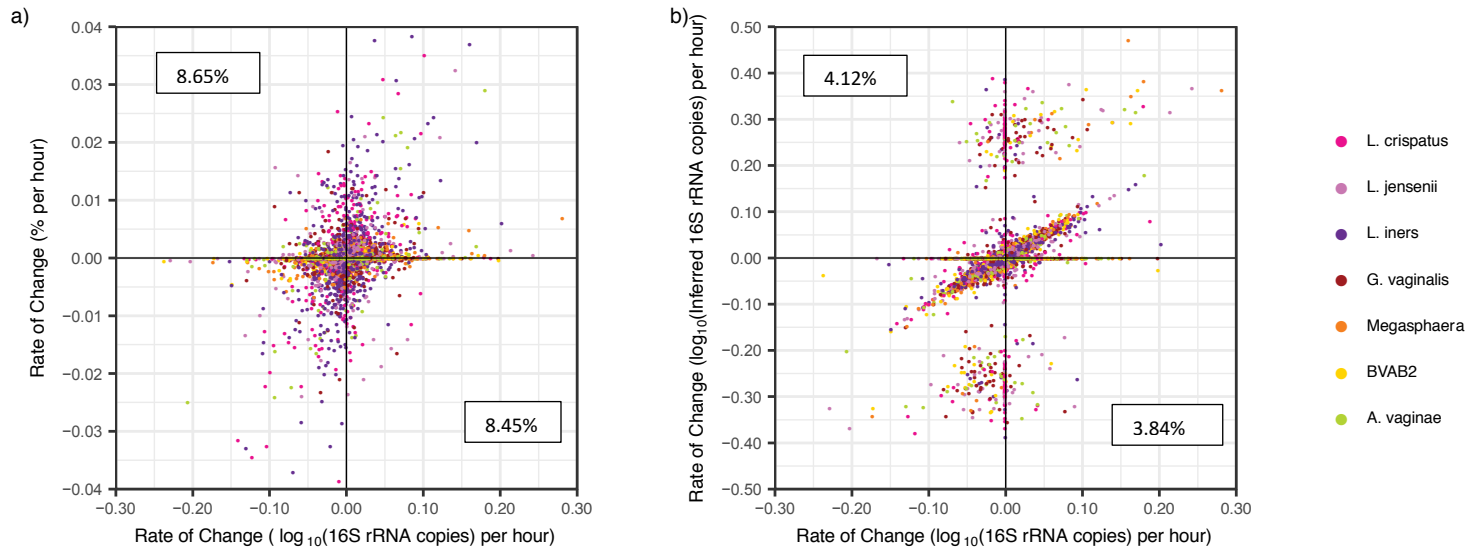
**Figure 6. Inferred concentrations allow somewhat accurate inference of kinetic changes between two sequential samples. a)** Scatter plot of change in relative abundance vs change absolute concentration shows poor correlation. Pearson correlation coefficient (pcc), r is 0.271 (P<2.2e-16). A high percentage of observed changes in relative abundance are in the opposite direction as those in absolute concentration (left upper and right lower error quadrants marked with percentages) **b)** Scatter plot of inferred concentration vs absolute concentration shows improved correlation. Both axes are plotted on a logarithmic scale. Pearson correlation coefficient (pcc), r is 0.392 (P<2.2e-16). Percentages correspond to the number of data points which fall within the error quadrants and are lower than for relative abundance. Inferred values misreport direction of kinetics less frequently.
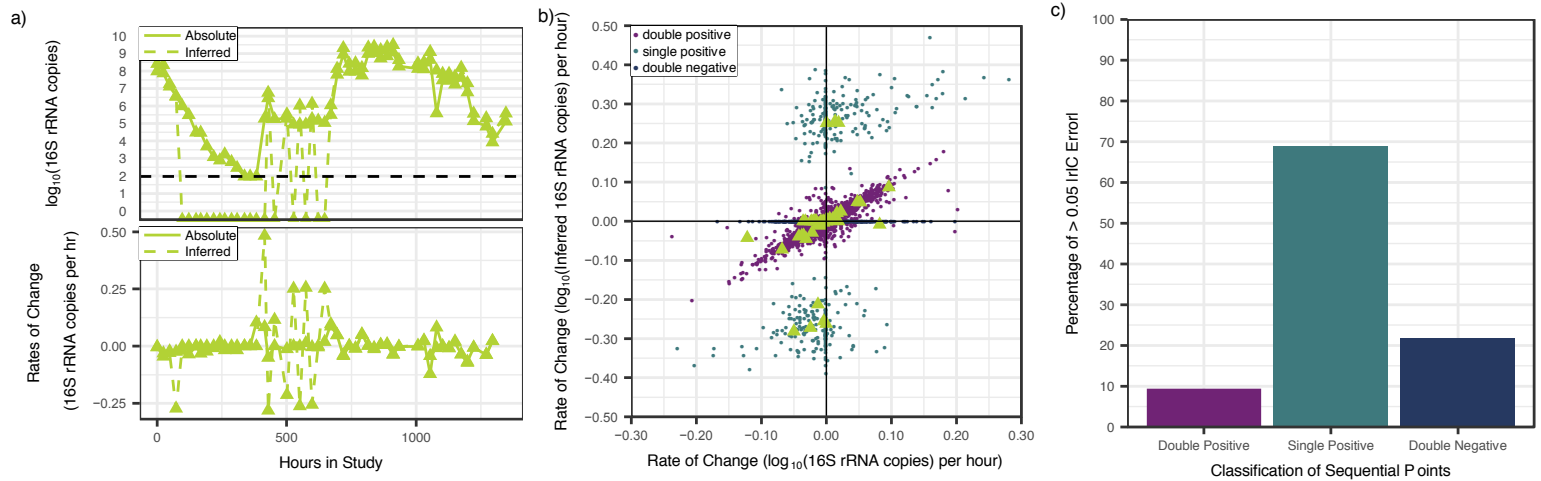
**Figure 7. Inferred concentration measures allow accurate inference of kinetic changes between two sequential non-negative samples. a)** Top: Levels of *A. vaginae* over time in a single participant (dotted is inferred and solid is absolute concentration); bottom: rate of change in levels of *A. vaginae* over time in the same participant (dotted is inferred and solid absolute concentration); divergence in swab to swab levels between inferred and absolute concentrations varies only when inferred concentration is zero in one of the sequential samples. **b)** Scatter plot of rate of change of inferred concentration as predicted by NGS vs qPCR observed values. Both axes are plotted on a logarithmic scale. Data is the same as in **Fig 6 a** and **b**. Triangles correspond to panel **c**. Points are colored according to whether consecutive samples were double positive (both >0 inferred concentration), single positive (one >0 and one 0 inferred concentration) or double negative (both 0 inferred concentration). Data points in which both samples are positive (no zeroes) are much more highly correlated (r is 0.876 , P<2.2e-16). **c)** A majority of rIC errors > 0.05 occur during transitions between positive and negative samples (one zero).