1  # On the robustness of inference of association with the gut microbiota in

2  # stool, swab and mucosal tissue samples

3

4  **Authors:**

5  Shan Sun[1], Xiangzhu Zhu[2], Xiang Huang[2], Harvey J. Murff[2], Reid M. Ness[2], Douglas L.
6  Seidner[3], Alicia Sorgen[1], Ivory Blakley[1], Chang Yu[4], Qi Dai[2], M. Andrea Azcarate-Peril[5],
7  Martha J. Shrubsole[2*], Anthony A. Fodor[1*]
8  [1]Department of Bioinformatics and Genomics, University of North Carolina at Charlotte,
9  Charlotte, NC, USA.
10  [2] Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA.
11  [3] Digestive Disease and Surgical Institute, Cleveland Clinic, Cleveland, OH 44195
12  [4]Department of Biostatistics, Vanderbilt University Medical Center, Nashville, TN, USA.
13  [5]Department of Medicine and Microbiome Core Facility, School of Medicine, University of
14  North Carolina, Chapel Hill, NC, USA.
15

16  **Corresponding Author:**

17  Anthony A. Fodor

18  Department of Bioinformatics and Genomics

19  University of North Carolina at Charlotte

20  9331 Robert D. Snyder Rd, Room 361

21  Charlotte, NC 28223

22  afodor@uncc.edu

23
24
25  *These authors contributed equally to this work.

26
27  **Key Words:** Gut microbiota, metagenome, stool, swab, mucosal tissue, host factors

28

29

## Abstract

The gut microbiota plays an important role in human health and disease. Stool, swab and mucosal tissue samples have been used in individual studies to survey the microbial community but the consequences of using these different sample types are not completely understood. We previously reported differences in microbial community composition with 16S rRNA amplicon sequencing between stool, swab and mucosal tissue samples. Here, we extended the previous study to a larger cohort and performed shotgun metagenome sequencing of 1,397 stool, swab and mucosal tissue samples from 240 participants. Consistent with previous results, taxonomic composition of stool and swab samples was distinct, but still more similar to each other than mucosal tissue samples, which had a substantially different community composition, characterized by a high relative abundance of the mucus metabolizers *Bacteroides* and *Subdoligranulum,* as well as bacteria with higher tolerance for oxidative stress such as *Escherichia*. As has been previously reported, functional profiles were more uniform across sample types than taxonomic profiles with differences between stool and swab samples smaller, but mucosal tissue samples remained distinct from the other two types. When the taxonomic and functional profiles of different sample types were used for inference in association with host phenotypes of age, sex, body mass index (BMI), antibiotics or non-steroidal anti-inflammatory drugs (NSAIDs) use, hypothesis testing using either stool or swab gave broadly similar results, but inference performed on mucosal tissue samples gave results that were generally less consistent with either stool or swab. Our study represents an important resource for the experimental design of studies aimed to understand microbiota perturbations specific to defined micro niches within the human intestinal tract.

**Introduction**

A growing number of studies have reported the essential roles of the human gut microbiota in human health and that microbiota alterations are associated with diseases including colorectal cancer, inflammatory bowel disease, obesity and diabetes [1-4]. The human colorectum is a complex system consisting of many microhabitats; studies have reported that the luminal and mucosal microbiota harbor heterogeneous microbial communities [5]. With the oxygen decline from the intestinal mucosa towards the lumen, anaerobic microorganisms are likely more abundant in luminal than mucosal environments [6]. On the other hand, the mucosal microbiota, directly adherent to the host tissue, may be more sensitive and respond more rapidly to localized changes in host tissues, compared to the luminal microbiota that is isolated from the loose mucus layer on the surface of the colorectal wall [7].

Stool samples are the most common biospecimen used to assess composition and functionality of the human gut microbiota in human research because of the large amount of biomass and the feasibility of collection; however, stool-derived profiles are more representative of luminal microorganisms than of mucosa-associated microbes. Mucosal tissue biopsy better characterizes mucosa-associated microbes but is less frequently used because of the invasive nature and accompanying risk of the procedure. Rectal swab may be used when stool samples are not practical to obtain, for example in the intensive care unit, and may collect a combination of both luminal and mucosal communities [8]. While stool and mucosal samples are generally distinct, there are mixed findings on the similarity between stool and swab samples [9-11]. Thus, different biospecimen types may be needed to sample microorganisms residing in different niches or to reflect different physiological conditions. For example, a study on colitis-induced inflammation in mouse

83 reported that microbial dysbiosis in the mucus layer was detected preceding colitis

84 while changes in stool microbiota were detected post-colitis [7].

85

86 Most of the studies assessing the variation of microbiota profile by biospecimen

87 type have focused on taxonomic composition characterized by 16S rRNA

88 amplicon sequencing. Previous literature of observed variation using shotgun

89 metagenomics is usually limited by the sample size, including our own previous

90 study [8]. Compared to the 16S rRNA amplicon sequencing, shotgun metagenome

91 sequencing utilizes total DNA instead of PCR products thus reducing the bias

92 introduced during processing. Moreover, metagenome sequencing not only

93 determines the taxonomic composition of the gut bacterial communities but also

94 generates information about functional information. With the increasing

95 application of shotgun metagenome sequencing in microbiota studies, a better

96 understanding of the metagenome variation across biospecimen types will help

97 investigators develop and interpret their experimental design.

98

99 In this study, we collected matched stool, rectal swab and mucosal tissue samples

100 from 240 study participants at two time points, which resulted in 1,397 shotgun

101 metagenomes. This is one of the largest studies comparing metagenomes of human

102 stool, rectal swab and colorectal mucosal tissue samples. We estimated the

103 biospecimen type variation of both metagenome taxonomy and functional

104 pathways. We also assessed whether the associations between taxa/pathways and

105 age, sex, body mass index (BMI), non-steroidal anti-inflammatory drugs (NSAIDs)

106 use and antibiotics use were consistent across the different sample types.

107

108 **Methods**

109 ***Study Population and Biospecimen Collection***

110    The Personalized Prevention of Colorectal Cancer Trial (PPCCT) was a double-

111    blind, placebo-controlled, randomized clinical trial designed to test the interaction

112    between a *TRPM7* genotype and reduction of the calcium/magnesium intake ratio

113    via magnesium supplementation on colorectal carcinogenesis biomarkers. Study

114    design and biospecimen collection have been previously described [8]. In brief,

115    participants were randomized to receive for 12 weeks either a personalized dose of

116    placebo (microcrystalline cellulose) or magnesium (magnesium glycinate).

117    Inclusion criteria included aged 40-85, personal history of colorectal polyps,

118    known *TRPM7* rs8042919 genotype, and daily intakes of calcium between 700-

119    2000 mg/day and the ratio of calcium to magnesium of 2.6 or greater. Exclusion

120    criteria included pregnancy, breastfeeding, use of medications that may interact

121    with magnesium, or personal history of cancer, colon resection or colectomy,

122    inflammatory bowel disease, organ transplantation, gastric bypass, chronic

123    diarrhea, chronic renal disease, hepatic cirrhosis, chronic ischemic heart disease, or

124    Type I diabetes. All study procedures were performed in accordance with relevant

125    guidelines and regulations as approved by the Vanderbilt Institutional Review

126    Board. The study is registered at ClinicalTrials.gov (NCT01105169).

127

128    Biospecimens were collected at home or in an in-person study visit at the

129    beginning of the trial (baseline) and at the conclusion of the study 12 weeks later

130    (mean $12.3 \pm 1.03$ weeks) [8]. Stool samples were collected by study participants

131    at home using a white plastic collection container covering the toilet bowl,

132    aliquoted by the participant into sterile cryovials, and stored in the home freezer

133    until transport with an ice pack to the study visit. Stool was collected up to 3 days

134    prior to the study visit. Rectal swabs and mucosal tissues were collected by the

135    study physician at the study visits. Rectal swabs were collected by inserting a

136    culturette swab through the anal canal, swabbing the distal rectal mucosa, and

137    placing the swab into a cryovial. Rectal mucosal samples were collected through

138    an anoscope using standard mucosal biopsy forceps and these samples were placed

139    into separate storage vials. All three biospecimen types were frozen at −80□°C

140    until use.

141

142    *DNA Isolation and Sequencing*

143    Samples were transferred to a 2 ml tube containing 200 mg of ≤106 μm glass beads

144    (Sigma, St. Louis, MO) and 0.3 ml of Qiagen ATL buffer (Valencia, CA),

145    supplemented with lysozyme (20 mg/ml) (Thermo Fisher Scientific, Grand Island,

146    NY). The suspension was incubated at 37°C for 1 h with occasional agitation.

147    Subsequently the suspension was supplemented with 600IU of proteinase K and

148    incubated at 60°C for 1 h. Finally, 0.3 ml of Qiagen AL buffer were added and a

149    final incubation at 70°C for 10 minutes was carried out. Bead beating was then

150    performed for 3 minutes in a Qiagen TissueLyser II at 30Hz. After a brief

151    centrifugation, supernatants were transferred to a new tube containing 0.3 ml of

152    ethanol. DNA was purified using a standard on-column purification method with

153    Qiagen buffers AW1 and AW2 as washing agents and eluted in 10mM Tris (pH

154    8.0).

155

156    Whole-genome shotgun metagenomics (WGS) DNA sequencing was performed as

157    previously described [8]. Briefly, 1 ng of genomic DNA was processed using the

158    Nextera XT DNA Sample Preparation Kit (Illumina). Next, fragmented and tagged

159    DNA was amplified using a limited-cycle PCR program. In this step index 1(i7)

160    and index 2(i5) were added between the downstream bPCR adaptor and the core

161    sequencing library adaptor, as well primer sequences required for cluster formation.

162    The DNA library was purified using Agencourt® AMPure® XP Reagent. Each

163    sample was quantified and normalized prior to pooling. The DNA library pool was

164　loaded on the Illumina platform reagent cartridge and on the Illumina HiSeq

165　instrument.

166

167　***Bioinformatics and Statistical Analyses***

168　Sequencing output from the Illumina HiSeq4000 platform was converted to fastq

169　format and demultiplexed using Illumina Bcl2Fastq 2.18.0.12. Quality control of

170　the demultiplexed sequencing reads was verified by FastQC. Human genome

171　contamination was removed from the shotgun metagenome sequencing reads with

172　KneadData. The number of reads before and after removing human genome

173　contamination is shown in Fig. S1. The taxonomic composition of the filtered

174　reads was characterized with MetaPhlAn2 [12] while the functional pathways were

175　annotated with HUMAnN2 against the UniRef database [13]. Unmapped reads

176　were excluded from the following analyses. PCoA ordination was generated with

177　Bray-Curtis dissimilarity based on genus composition and functional pathway

178　abundance respectively with function 'capscale' in the R package 'vegan'. The

179　PERMANOVA test was performed with the function 'adonis' in the same package.

180　For each individual genus or pathway, we built linear mixed effects models with

181　the function 'lme' in R package 'nlme' with the aim of examining differences

182　between the modes based on sample type variation. The genera and pathways with

183　presence <10% in all samples were excluded to avoid spurious results and P-values

184　were adjusted with the Benjamini-Hochberg method for multiple testing.

185

186　Model 1 was used to test the associations between the metagenome and

187　biospecimen types (stool, swab or mucosal tissue). Model 1 was performed for

188　each pair of sample types to get the direction of changes and adjusted for host

189　factors.

190

$$\text{Genus/pathway} = \text{sample\_type}+\text{treatment*time\_point}+\text{antibiotics use} \quad (1)$$
$$+\text{age}+\text{sex}+\text{BMI}+\text{NSAIDs use}+(1/\text{participant})$$

191

192 In this model, sample type, treatment, time point, age, sex, BMI, antibiotics and

193 NSAIDs use were fixed effects while participant ID was a random effect. Using

194 pairwise models allowed for direct comparison between sample types. The

195 significance was determined as <10% FDRs corrected with Benjamini-Hochberg

196 method. Significant genera and pathways identified in this model were plotted as

197 heatmaps with the function 'pheatmap'.

198

199 Model 2 was used to test the associations between metagenome and host factors in

200 each sample type.

$$\text{Genus/pathway} = \text{treatment*time\_point}+\text{antibiotics use} \quad (2)$$
$$+\text{age}+\text{sex}+\text{BMI}+\text{NSAIDs use} +(1/\text{participant})$$

201 In this model, treatment, time point, age, sex, BMI, antibiotics and NSAIDs use

202 were fixed effects while participant ID is a random effect. The correlations

203 between inferences (-log10(P)) produced in different sample types were tested with

204 Spearman correlations and the plots were generated with 'ggplot2'.

205

206 Because of the compositional nature of the shotgun metagenome sequencing data,

207 we also utilized ALDEx2 [14] which uses Bayesian methods and a geometric mean

208 based normalization to minimize compositional artifacts to confirm our

209 observations. Because ALDEx2 does not support models adjusted for covariates,

210 the associations were tested with one variable models.

211

212 **Results**

213 ***Taxonomic composition of metagenomes was associated with sample types***

214  After quality control, there were 1,397 stool, swab and mucosal tissue

215  metagenomes from 240 participants. We characterized the taxonomic composition

216  and functional pathways of the metagenomes and found substantial variation by

217  sample type. Shannon diversity at the genus level was significantly different

218  between sample types, with mucosal tissue samples of the lowest diversity and

219  swab the highest (Fig. 1a). PCoA ordinations of genus composition showed a

220  distinct cluster of mucosal tissue samples (Fig. 1b). A PCoA ordination in which

221  mucosal tissue samples were excluded in order to better visualize the stool and

222  swab samples showed clear separation as well (Fig. 1c). A PERMANOVA test

223  indicated that the genus composition was significantly associated with sample type

224  (P=0.001, with 999 permutations). The differences across stool, swab and mucosal

225  tissue samples explained 31.6% of the variance, while the differences between

226  stool and swab explained 5%, further supporting the observation that mucosal

227  tissue samples were more distinct compared to stool and swab. Microbial

228  taxonomic composition at other levels from phylum to species levels were also

229  significantly associated with sample type (Table S1). The PERMANOVA tests and

230  PCoA ordinations demonstrate that the microbial metagenomes sampled with

231  different methods were different at the community level.

232

233  In order to identify differentially abundant taxa, we used a linear mixed-effects

234  models to compare the sample types in pairs (Model 1). Among the 60 genera with

235  presence in >10% samples, 56 were different between at least one pair of samples,

236  with 35 significantly different between stool and swab samples, 53 between stool

237  and tissue, and 51 between swab and tissue (Fig. 2). Because the sequencing

238  depths were different between sample types (Fig. S1), we also utilized an analysis

239  pipeline based on ALDEx2, which attempts to explicitly correct for compositional

240  artifacts. The differential abundance of the 56 taxa across sample types were

241    supported by results from ALDEx2, except for *Paraprevotella* and an unknown

242    genus of the Clostridiaceae family (Table S2). P-values from the two methods were

243    generally consistent (Fig. S2a). Tissue samples had higher relative abundance of

244    *Bacteroides*, *Subdoligranulum*, *Escherichia*, *Blautia* and unclassified genera of the

245    families *Propionibacteriaceae* and *Acidaminococcaceae*. Compared to stool

246    samples, swab samples were enriched in *Propionibacterium*, *Campylobacter*,

247    *Porphyromonas*, *Prevotella*, *Clostridium*, *Streptococcus* and had lower abundance

248    of *Methanobrevibacter*, *Dialister*, *Adlercreutzia*, *Haemophilus*, *Klebsiella*,

249    *Akkermansia*, *Alistipes* and *Paraprevotella*.

250

### *Functional pathways of metagenomes were associated with sample types*

252    The metagenomes of mucosal tissue samples had a higher number of reads that

253    could not be mapped to the UniRef databases after removing host sequences (54%

254    compared to 30% for stool and 28% for swab samples), indicating that the mucosal

255    tissue microbiota was less represented in the current database. The number of

256    microbial pathways was lower in mucosal tissues compared to other samples (Fig.

257    3a). The PCoA ordinations of functional pathways showed a similar specific

258    cluster of mucosal tissue samples (Fig. 3b), while the stool and swab samples were

259    less separated compared to the PCoA ordination based on genus composition

260    (Fig.3c). A PERMANOVA test indicated that functional pathways were also

261    significantly different across sample types (stool, swab and mucosal tissue: $R^2 =$

262    0.273, P=0.001; stool and swab: $R^2 = 0.048$, P=0.001). We again used a linear

263    mixed effects model to identify the differential functional pathways between

264    samples. In 343 functional pathways with presence in >10% samples, 318 were

265    significantly different between at least one pair of samples, with 269 of differential

266    abundance for stool-swab comparison, 222 for stool-tissue and 233 for swab-tissue

267    (Fig. 4). Among the 318 significant pathways, only 8 were not supported by the

268    analysis of ALDEx2 (Table S3; Fig. S2b).

269

270    ***The impact of sample type on the associations between the taxonomic and***

271    ***functional profiles and host factors***

272    We built separate models in each sample type to estimate whether the associations

273    with taxonomic composition were consistent across sample types for the host

274    factors age, sex, BMI, antibiotics use and NSAIDs use. The associations between

275    genera and host factors were very highly correlated between stool and swab

276    samples (Fig. 5: left panels) with Spearman's correlation coefficients of p-value vs.

277    p-value ranging from 0.501 for BMI to 0.75 for sex. The associations between

278    stool and mucosal tissue samples (Fig. 5: middle panels) were significantly

279    correlated except for sex with a P-value cutoff of 0.05, while the associations

280    between swab and mucosal tissue samples (Fig. 5: right panels) were significantly

281    correlated for BMI, antibiotics use and NSAIDs use but not for age or sex.

282

283    The same models were used for analyzing the robustness of the associations

284    between pathways and host factors (Fig. 6). As was the case for taxa, the

285    associations between pathways and host factors observed in stool and swab sample

286    types were all highly positively correlated (Fig 6: left panels). However,

287    comparisons between mucosal tissue and stool (Fig 6: middle panels) and swab

288    (Fig 6: right panels) samples showed that the correlations were less consistent,

289    including positive correlation with a smaller coefficient, negative correlation and

290    no correlation. These observations were generally consistent when using ALDEx2

291    for statistical modeling instead of the linear models for both taxonomic

292 composition and functional pathways that inference with stool and swab are more

293 consistent than with mucosal tissue (Table S4 and S5).

294

## 295 Discussion

296 A better understanding of the associations between the human gut microbiome and

297 disease is essential for developing potential early detection and intervention

298 methods utilizing the microbiome. The stool, rectal swab and mucosal tissue

299 biospecimen types we examined in this study sample microhabitats in which

300 different microbial communities reside. With 1397 matched stool, rectal swab and

301 mucosal tissue metagenomes for 240 participants, our dataset provided a great

302 opportunity for analyzing the variations of these three matched biospecimens from

303 the same participants. Unsurprisingly, we found that microbial taxonomic

304 composition and functional pathways were different across the three biospecimen

305 types, with the mucosal tissue metagenome more distinct from stool and swab. In

306 general, the inference of host factor and microbiome associations were highly

307 consistent between stool and rectal swab but not for mucosal tissue.

308

309 The mucosal tissue microbiome had lower alpha diversity and low abundance of

310 most microbes, but was enriched in *Bacteroides, Subdoligranulum, Escherichia*

311 and *Propionibacteriaceae*. *Bacteroides thetaiotaomicron, B. caccae, B. fragilis* and

312 *B. vulgatus* are well known mucin degraders and rely on mucin and other host-

313 derived glycans for colonization [15]. *Propionibacterium* (phylum Actinobacteria)

314 and *Escherichia* (phylum Proteobacteria) were higher in mucosal tissue and swab

315 compared to stool samples, which could be explained by their higher oxygen

316 tolerance. The enrichment of Actinobacteria and Proteobacteria in the mucosa-

317 associated microbiota has been reported in correlation with the intestinal radial

318 colonic oxygen gradient that influences microbiota composition based on their

319  ability to tolerate the oxidative stress [16]. The higher alpha diversity in the rectal

320  swab microbiome compared to the stool and mucosal tissue microbiome is

321  consistent with our previous study [8] and could be explained by swab sampling

322  from both luminal and mucosal microbes [9].

323

324  Similar to taxonomic composition, the functional pathways in stool and rectal swab

325  samples were more similar to each other than mucosal tissue samples. The number

326  of sequencing reads from the mucosal tissue was smaller compared to stool and

327  rectal swab samples due to lower microbial biomass and a higher percentage of

328  human genome DNA contamination (Fig. S1). This could contribute to the

329  observed lower taxonomic and functional diversity in mucosal tissue microbiome

330  compared to stool and rectal swab samples. Compositional artifacts associated with

331  reduced sequencing depth may therefore explain some of the differences we

332  observed between mucosal tissue and stool and swab samples. These differences

333  did persist even when using the compositionally aware pipeline ALDEx2, but no

334  statistical approach can perfectly compensate for large differences in sequencing

335  depth. ALDEx2 does not allow for inclusion of covariates or adjusting for random

336  effects from the same subject and that might explain the differences between the

337  ALDEx2 and linear models. Future research will be needed to explore how much

338  of the differences between mucosal tissues and stool and swab in both community

339  and gene composition and inference can be explained by these compositional

340  differences.

341

342  Stool, swab and mucosal microbiota were enriched for different pathways,

343  reflecting the niche adaption of different microbial communities. Mucosal

344  microbiota was relatively enriched for pathways related to glycolysis and

345  biosynthesis pathways involved in the generation of amino acid L–isoleucine,

346  nucleosides adenosine, guanosine and inosine, and fatty acids gondoate and *cis*-

347  vaccenate (one of the major unsaturated fatty acids, responsible for membrane

348  phospholipid homeostasis in bacteria[17]). The stool and rectal swab microbiomes

349  differed in the pathway related to peptidoglycan, CDP–diacylglycerol,

350  UDP–N–acetylmuramoyl–pentapeptide, galactose, stachyose, L–arginine, purine

351  and pyrimidine. Because a large number of functional genes remained unexplored,

352  future expansion of database could provide a better knowledge of the functional

353  differences between these sample types.

354  In order to determine whether the biospecimen type influence the inference of

355  associations between the gut microbiome and host factors, we analyzed microbial

356  associations with age, sex, BMI, antibiotics and NSAIDs use in each of the three

357  sample types. We found that inferences performed with stool and rectal swab

358  samples were highly correlated with each other for both taxonomic composition

359  and functional pathways, while inference with mucosal tissue was more distinct

360  especially for functional pathways. The relatively poor consistency between the

361  mucosal tissue microbiome and the stool and rectal swab microbiome potentially

362  reflects the niche differences that affect microbial interactions with the

363  environment. It is also possible that the mucus barrier between the mucosal tissue

364  and the lumen makes the mucosal tissue microbiome more sensitive to some host

365  changes that were reflected in the mucosal tissues. For example, a previous study

366  reported that the excessive secretion of mucus glycan could lead to the increase of

367  *Akkermansia* and *Bacteroides* abundance in mucosal tissue but was only extended

368  to stool with an altered mucus barrier [7]. As is the case for comparisons of relative

369  abundance, models of inference are also sensitive to compositional artifacts

370     associated with sequencing depth, although in our study comparisons based on

371     ALDEx2 yielded broadly similar results to comparisons based on compositionally

372     naïve mixed linear models.

373

374     We note that this study was conducted in individuals with a history of colorectal

375     polyps, so the conclusions may not be generalizable to individuals without a

376     history of polyps. However, all the participants were polyp-free when

377     biospecimens were collected. Our work represents the largest study to date to

378     explicitly compare these sample types and should provide a useful guide to

379     investigators in the design and interpretation of human studies of the gut

380     microbiota.

381

## 382     **Conclusion**

383     Our study shows that the stool, swab and mucosal tissue microbiota are of different

384     taxonomic and functional profiles, but the stool and swab microbiota are generally

385     more similar compared to that of mucosal tissue. When analyzing the associations

386     between microbiota and host factors of age, sex, BMI, antibiotics or NSAIDs use

387     in each sample type, the inference on stool and swab samples were also more

388     consistent than the inference on mucosal samples. Our study suggests that not only

389     the taxonomic and functional profiles varied by sample types but the inference on

390     their associations with host factors were depending on the sample type as well.

391

## Author Contributions

410     QD, CY, MJS and AAF contributed to study conception, design, and supervision. XZ, HJM,
411     RMN, DLS, MAAP, QD and MJS contributed to acquisition of data. XZ and MJS provided
412     administrative, technical, or material support. SS, XZ, HX, AS, IB, CY, DQ, MJS and AAF
413     contributed to analysis and interpretation of data. All authors contributed to writing, review,
414     and/or revision of the manuscript and approved the final manuscript.
415
416
417

## Competing interests

418     The authors declare no competing interests.
419

## Data sharing statement

420     The metagenomes sequences analyzed in this study are available at NBCI with accession ID
421     PRJNA693850. Scripts used in this study are available at
422     https://github.com/ssun6/StoolSwabTissue.
423

## REFERENCES

427   1.     Arthur JC, Gharaibeh RZ, Mühlbauer M, Perez-Chanona E, Uronis JM, McCafferty J,
428         Fodor AA, Jobin C: Microbial genomic analysis reveals the essential role of
429         inflammation in bacteria-induced colorectal cancer. *Nat Commun* 2014, 5**:**4724.
430   2.     Kostic AD, Xavier RJ, Gevers D: The microbiome in inflammatory bowel disease:
431         current status and the future ahead. *Gastroenterology* 2014, 146**:**1489-1499.
432   3.     Graham C, Mullen A, Whelan K: Obesity and the gastrointestinal microbiota: a review of
433         associations and mechanisms. *Nutrition reviews* 2015, 73**:**376-385.
434   4.     Qin J, Li Y, Cai Z, Li S, Zhu J, Zhang F, Liang S, Zhang W, Guan Y, Shen D: A
435         metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature* 2012,
436         490**:**55.
437   5.     Donaldson GP, Lee SM, Mazmanian SK: Gut biogeography of the bacterial microbiota.
438         *Nature Reviews Microbiology* 2016, 14**:**20-32.
439   6.     Espey MG: Role of oxygen gradients in shaping redox relationships between the human
440         intestine and its microbiota. *Free Radical Biology and Medicine* 2013, 55**:**130-140.
441   7.     Glymenaki M, Singh G, Brass A, Warhurst G, McBain AJ, Else KJ, Cruickshank SM:
442         Compositional changes in the gut mucus microbiota precede the onset of colitis-induced
443         inflammation. *Inflammatory bowel diseases* 2017, 23**:**912-922.

444 8.  Jones RB, Zhu X, Moan E, Murff HJ, Ness RM, Seidner DL, Sun S, Yu C, Dai Q, Fodor
445     AA: Inter-niche and inter-individual variation in gut microbial community assessment
446     using stool, rectal swab, and mucosal samples. *Scientific reports* 2018, 8**:**4139.
447 9.  Choudhury R, Kleerebezem M, Middelkoop A, Bolhuis JE: Legitimate and reliable
448     determination of the age-related intestinal microbiome in young piglets; rectal swabs and
449     fecal samples provide comparable insights. *Frontiers in Microbiology* 2019, 10**:**1886.
450 10. Bassis CM, Moore NM, Lolans K, Seekatz AM, Weinstein RA, Young VB, Hayden MK:
451     Comparison of stool versus rectal swab samples and storage conditions on bacterial
452     community profiles. *BMC microbiology* 2017, 17**:**1-7.
453 11. Fair K, Dunlap DG, Fitch A, Bogdanovich T, Methé B, Morris A, McVerry BJ, Kitsios
454     GD: Rectal swabs from critically ill patients provide discordant representations of the gut
455     microbiome compared to stool samples. *Msphere* 2019, 4**:**e00358-00319.
456 12. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A,
457     Huttenhower C, Segata N: MetaPhlAn2 for enhanced metagenomic taxonomic profiling.
458     *Nature methods* 2015, 12**:**902-903.
459 13. Franzosa EA, McIver LJ, Rahnavard G, Thompson LR, Schirmer M, Weingart G, Lipson
460     KS, Knight R, Caporaso JG, Segata N: Species-level functional profiling of metagenomes
461     and metatranscriptomes. *Nature methods* 2018, 15**:**962.
462 14. Fernandes AD, Macklaim JM, Linn TG, Reid G, Gloor GB: ANOVA-like differential
463     expression (ALDEx) analysis for mixed population RNA-Seq. *PLoS One* 2013, 8**:**e67019.
464 15. Martens EC, Chiang HC, Gordon JI: Mucosal glycan foraging enhances fitness and
465     transmission of a saccharolytic human gut bacterial symbiont. *Cell host & microbe* 2008,
466     4**:**447-457.
467 16. Albenberg L, Esipova TV, Judge CP, Bittinger K, Chen J, Laughlin A, Grunberg S,
468     Baldassano RN, Lewis JD, Li H: Correlation between intraluminal oxygen gradient and
469     radial partitioning of intestinal microbiota. *Gastroenterology* 2014, 147**:**1055-1063.
470     e1058.
471 17. Zhang Y-M, Rock CO: Membrane lipid homeostasis in bacteria. *Nature Reviews*
472     *Microbiology* 2008, 6**:**222-233.
473
474

475 **Figure legends**

476 Fig. 1. Alpha-diversity and PCoA ordinations of the taxonomic composition of

477 microbial metagenomes at the genus level. Color indicates the sample types. (a)

478 Alpha diversity across sample types. Differences between sample types were tested

479 with Wilcoxon's test. (b) Mucosal tissue samples formed a distinct cluster from

480 stool and swab samples. (c) Separation of stool and swab samples.

481

482    Fig. 2. Heatmap of genera that were significantly different between sample types

483    (FDR<0.05). Keys indicate the z-scores of averaged taxonomic abundance.

484

485    Fig. 3. The number of pathways and PCoA ordinations of functional pathways of

486    microbial metagenomes. Color indicates the sample types. (a) The number of

487    pathways across samples. (b) Mucosal tissue samples formed a distinct cluster

488    from stool and swab samples. (c) visualization of only stool and swab samples.

489

490    Fig. 4. Heatmap of functional pathways that were significantly different between

491    sample types (FDR<0.05). Keys indicate z-scores of averaged abundance.

492

493    Fig. 5. Correlations between the genus composition inference for age (a), sex (b),

494    BMI (c), antibiotics use (d) and NSAIDs use (e) between pairwise sample types.

495    The axes were the –log10 transformation of p-values from the model 2 described in

496    methods

497

498    Fig 6. Correlations between the functional pathways inference for age (a), sex (b),

499    BMI (c), antibiotics use (d), NSAIDs use (e) and between pairwise sample types.

500    The axes were the –log10 transformation of p-values from the model 2 described in

501    methods.

502

503

Fig. 1. Alpha-diversity and PCoA ordinations of taxonomic composition of microbial metagenomes at the genus level. Color indicates the sample types. (a) Alpha diversity across sample types. Differences between sample types were tested with Wilcoxon's test. (b) Mucosal tissue samples formed a distinct cluster from stool and swab samples. (c) Separation of stool and swab samples.

Fig. 2. Heatmap of genera that were significantly different between sample types (FDR<0.05). Keys indicate z-scores of averaged taxonomic abundance.

Fig. 3. The number of pathways and PCoA ordinations of functional pathways of microbial metagenomes. Color indicates the sample types. (a) The number of pathways across samples. (b) Mucosal tissue samples formed a distinct cluster from stool and swab samples. (c) visualization of only stool and swab samples.

PWY-7184: pyrimidine deoxyribonucleotides de novo biosynthesis I
PWY-6121: 5–aminoimidazole ribonucleotide biosynthesis I
PWY-6122: 5–aminoimidazole ribonucleotide biosynthesis II
PWY-6277: superpathway of 5–aminoimidazole ribonucleotide biosynthesis
PWY-3001: superpathway of L–isoleucine biosynthesis I
PWY-5686: UMP biosynthesis
PWY-841: superpathway of purine nucleotides de novo biosynthesis I
CALVIN-PWY: Calvin–Benson–Bassham cycle
PWY0-166: superpathway of pyrimidine deoxyribonucleotides de novo biosynthesis (E. coli)
PWY-7229: superpathway of adenosine nucleotides de novo biosynthesis I
THRESYN-PWY: superpathway of L–threonine biosynthesis
PWY-7208: superpathway of pyrimidine nucleobases salvage
COA-PWY-1: coenzyme A biosynthesis II (mammalian)
PEPTIDOGLYCANSYN-PWY: peptidoglycan biosynthesis I (meso–diaminopimelate containing)
PWY-6387: UDP–N–acetylmuramoyl–pentapeptide biosynthesis I (meso–diaminopimelate containing)
PWY-7197: pyrimidine deoxyribonucleotide phosphorylation
PWY-4242: pantothenate and coenzyme A biosynthesis III
PWY-6386: UDP–N–acetylmuramoyl–pentapeptide biosynthesis II (lysine–containing)
PWY-5667: CDP–diacylglycerol biosynthesis I
PWY0-1319: CDP–diacylglycerol biosynthesis II
PWY0-162: superpathway of pyrimidine ribonucleotides de novo biosynthesis
ILEUSYN-PWY: L–isoleucine biosynthesis I (from threonine)
NONMEVIPP-PWY: methylerythritol phosphate pathway I
PWY-724: superpathway of L–lysine, L–threonine and L–methionine biosynthesis II
PWY-1042: glycolysis IV (plant cytosol)
PWY-5101: L–isoleucine biosynthesis II
PWY-7663: gondoate biosynthesis (anaerobic)
PWY-7220: adenosine deoxyribonucleotides de novo biosynthesis II
PWY-7222: guanosine deoxyribonucleotides de novo biosynthesis II
PWY-6151: S–adenosyl–L–methionine cycle I
SO4ASSIM-PWY: sulfate reduction I (assimilatory)
PWY-5973: cis–vaccenate biosynthesis
GLYCOLYSIS: glycolysis I (from glucose 6–phosphate)
PWY-5484: glycolysis II (from fructose 6–phosphate)
ANAGLYCOLYSIS-PWY: glycolysis III (from glucose)
PWY-6123: inosine–5'–phosphate biosynthesis I
PWY0-1586: peptidoglycan maturation (meso–diaminopimelate containing)
PWY-5695: urate biosynthesis/inosine 5'–phosphate degradation
PWY-6124: inosine–5'–phosphate biosynthesis II
DTDPRHAMSYN-PWY: dTDP–L–rhamnose biosynthesis I
RHAMCAT-PWY: L–rhamnose degradation I
VALSYN-PWY: L–valine biosynthesis
PWY-7111: pyruvate fermentation to isobutanol (engineered)

PWY-5125: superpathway of guanosine nucleotides de novo biosynthesis II

1CMET2-PWY:

PWY-6126: superpathway of adenosine nucleotides de novo biosynthesis II
PWY-7228: superpathway of guanosine nucleotides de novo biosynthesis I
BRANCHED-CHAIN-AA-SYN-PWY: superpathway of branched amino acid biosynthesis
PWY0-845: superpathway of pyridoxal 5'–phosphate biosynthesis and salvage
PANTOSYN-PWY: pantothenate and coenzyme A biosynthesis I
PWY-5097: L–lysine biosynthesis VI
PWY-6703: preQ0 biosynthesis
ARGININE-SYN4-PWY: L–ornithine de novo  biosynthesis
PWY-5659: GDP–mannose biosynthesis
PWY-6385: peptidoglycan biosynthesis III (mycobacteria)
COA-PWY: coenzyme A biosynthesis I
PANTO-PWY: phosphopantothenate biosynthesis I
PWY-7219: adenosine ribonucleotides de novo biosynthesis
PWY-2942: L–lysine biosynthesis III
PYRIDOXSYN-PWY: pyridoxal 5'–phosphate biosynthesis I
NONOXIPENT-PWY: pentose phosphate pathway (non–oxidative branch)
PWY-6700: queuosine biosynthesis
PWY-1269: CMP–3–deoxy–D–manno–octulosonate biosynthesis I
PWY-3841: folate transformations II
HISTSYN-PWY: L–histidine biosynthesis
GLCMANNANAUT-PWY: superpathway of N–acetylglucosamine, N–acetylmannosamine and N–acetylneuraminate degradation
PHOSLIPSYN-PWY: superpathway of phospholipid biosynthesis I (bacteria)
GALACTUROCAT-PWY: D–galacturonate degradation I
PWY66-400: glycolysis VI (metazoan)
ASPASN-PWY: superpathway of L–aspartate and L–asparagine biosynthesis
GLUCUROCAT-PWY: superpathway of &beta;–D–glucuronide and D–glucuronate degradation
PWY-6168: flavin biosynthesis III (fungi)
PWY-7242: D–fructuronate degradation
GALACT-GLUCUROCAT-PWY: superpathway of hexuronide and hexuronate degradation
PWY-5103: L–isoleucine biosynthesis III
THISYNARA-PWY: superpathway of thiamin diphosphate biosynthesis III (eukaryotes)
PWY-6609: adenine and adenosine salvage III
SER-GLYSYN-PWY: superpathway of L–serine and glycine biosynthesis I
TRPSYN-PWY: L–tryptophan biosynthesis
ANAEROFRUCAT-PWY: homolactic fermentation
PYRIDNUCSYN-PWY: NAD biosynthesis I (from aspartate)
MET-SAM-PWY: superpathway of S–adenosyl–L–methionine biosynthesis
PWY-6270: isoprene biosynthesis I
PWY-5188: tetrapyrrole biosynthesis I (from glutamate)
COMPLETE-ARO-PWY: superpathway of aromatic amino acid biosynthesis
DENOVOPURINE2-PWY: superpathway of purine nucleotides de novo biosynthesis II
METSYN-PWY: L–homoserine and L–methionine biosynthesis
PWY-6737: starch degradation V
PWY-7187: pyrimidine deoxyribonucleotides de novo biosynthesis II
ARO-PWY: chorismate biosynthesis I
PWY-6163: chorismate biosynthesis from 3–dehydroquinate
TRNA-CHARGING-PWY: tRNA charging
PWY-6897: thiamin salvage II
PWY66-422: D–galactose degradation V (Leloir pathway)
PWY-6527: stachyose degradation
PWY-7400: L–arginine biosynthesis IV (archaebacteria)
ARGSYN-PWY: L–arginine biosynthesis I (via L–ornithine)
RIBOSYN2-PWY: flavin biosynthesis I (bacteria and plants)
PWY-6317: galactose degradation I (Leloir pathway)
PWY-7357: thiamin formation from pyrithiamine and oxythiamine (yeast)
PWY-7282: 4–amino–2–methyl–5–phosphomethylpyrimidine biosynthesis (yeast)
PWY-7560: methylerythritol phosphate pathway II
PWY0-1296: purine ribonucleosides degradation

tissue    stool    swab

Fig. 5. Correlations between the genus composition inference for age (a), sex (b), BMI (c), antibiotics use (d) and NSAIDs use (e) between pairwise sample types. The axes were the −log10 transformation of p-values from the model2 described in methods

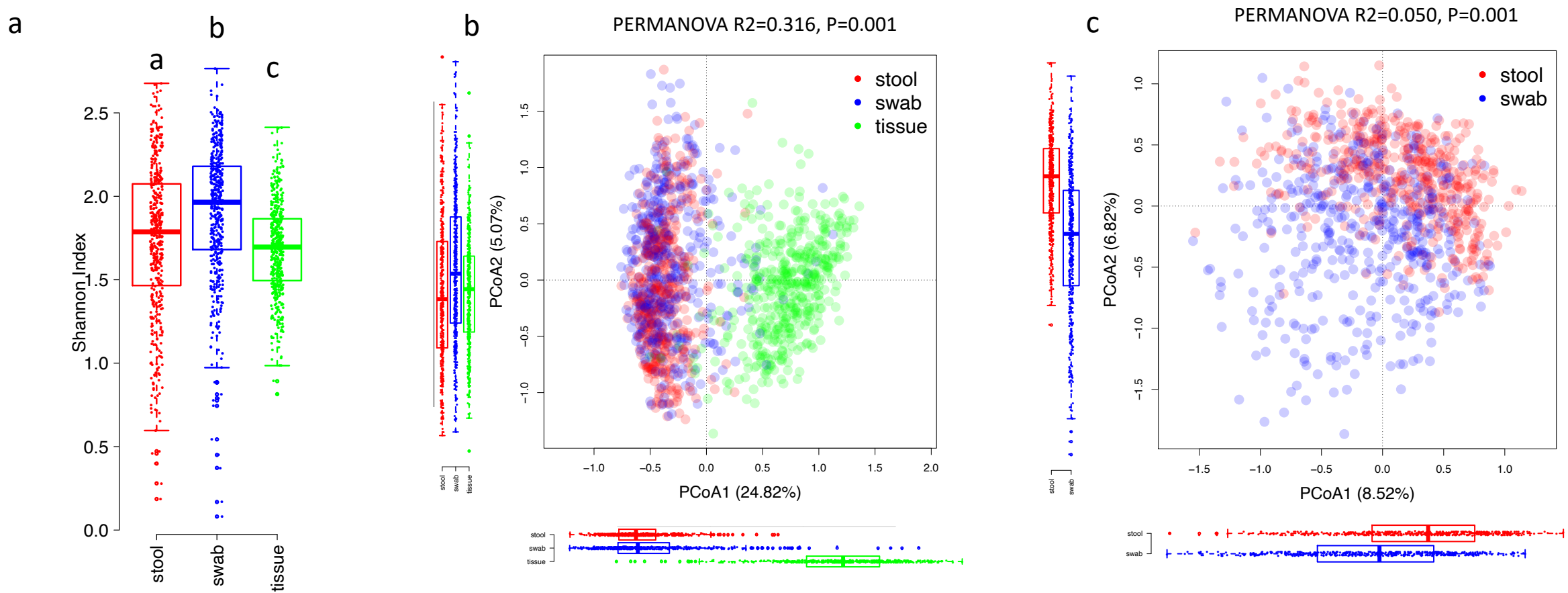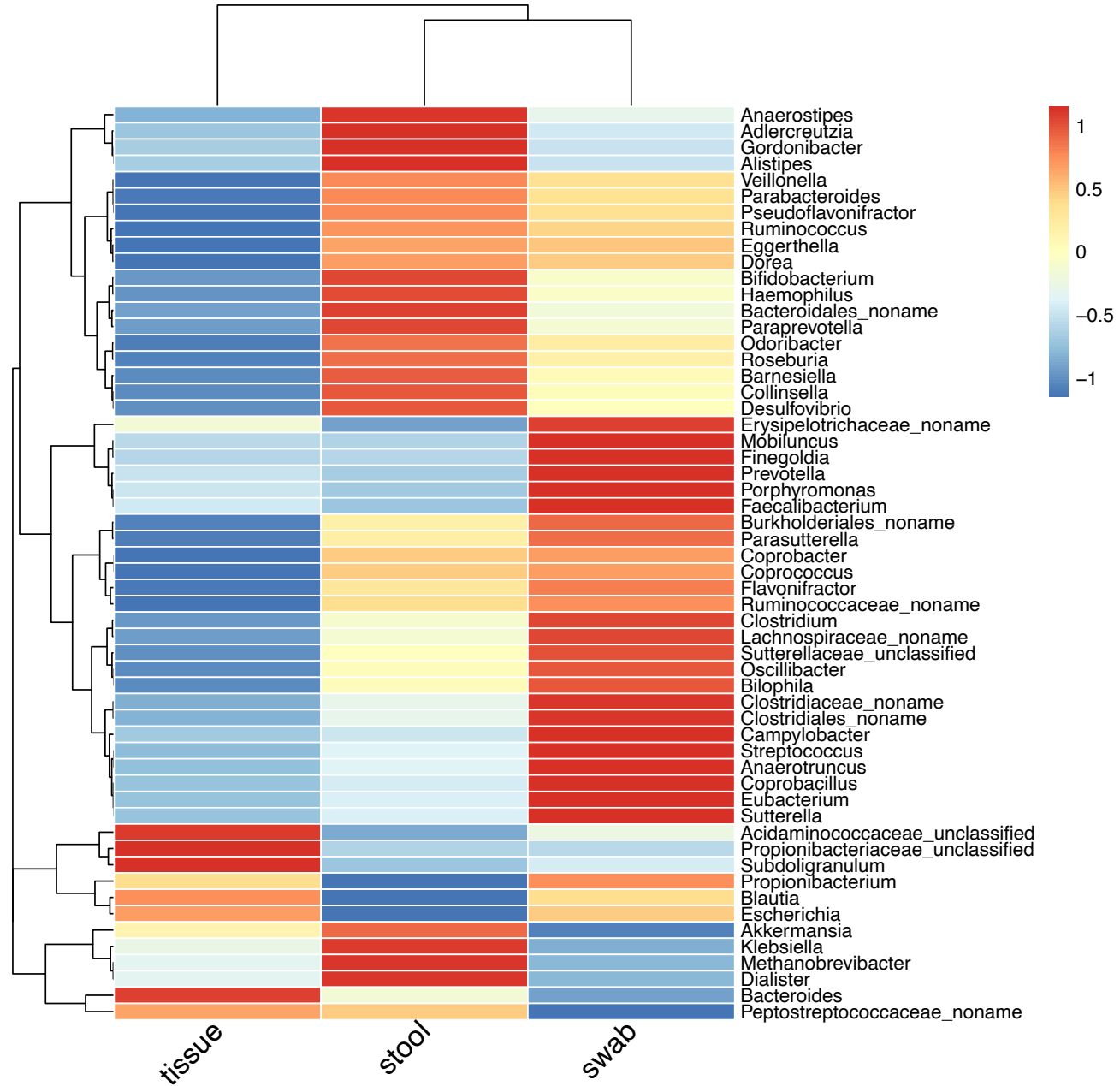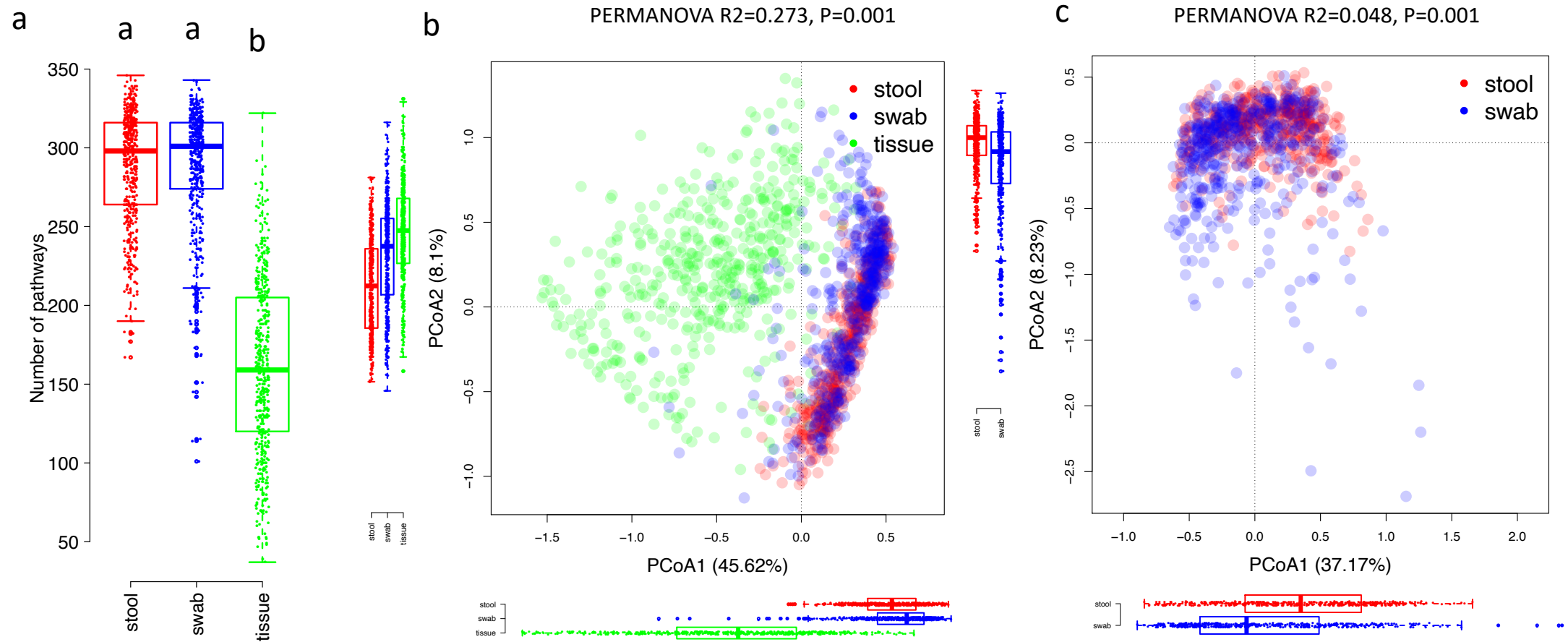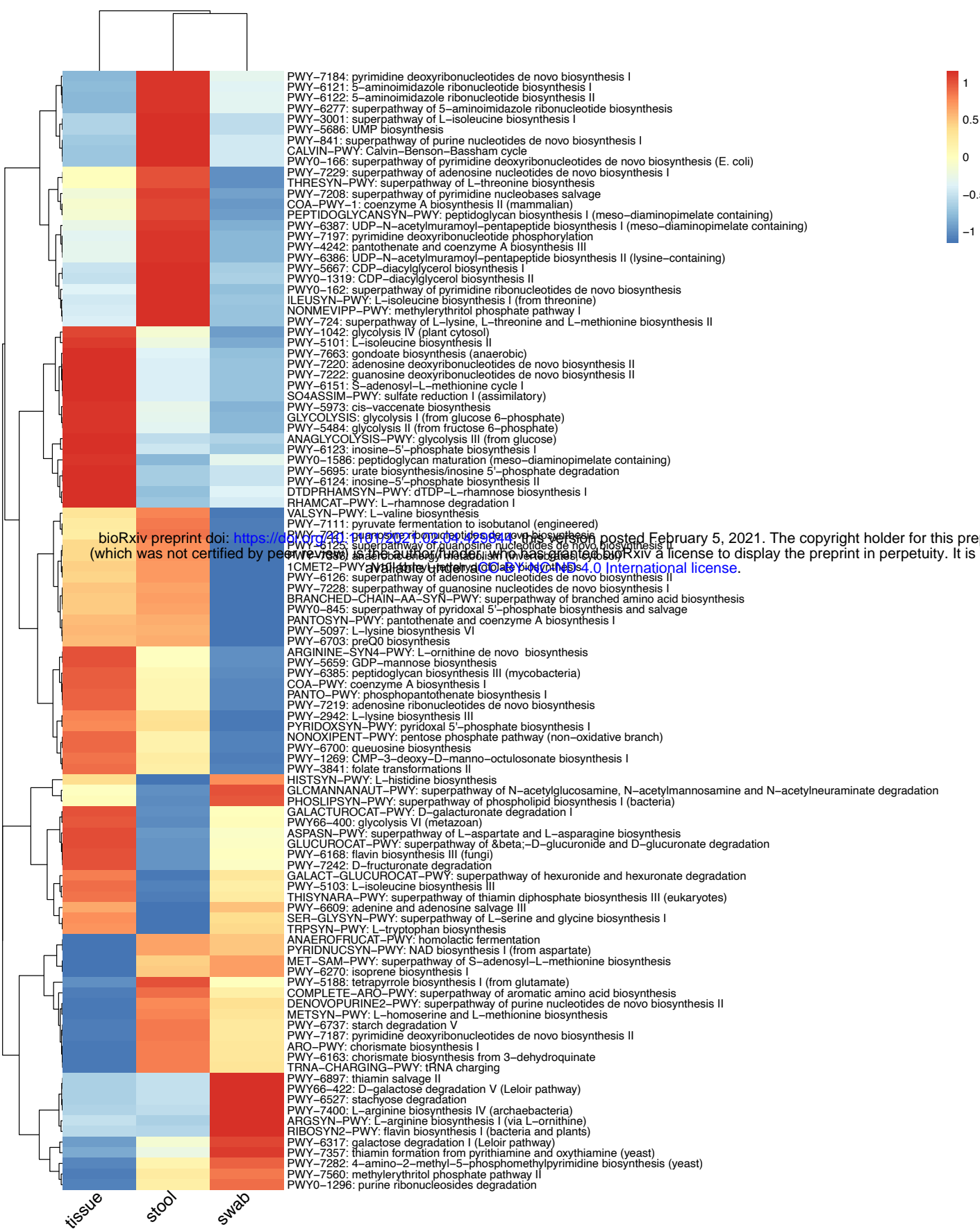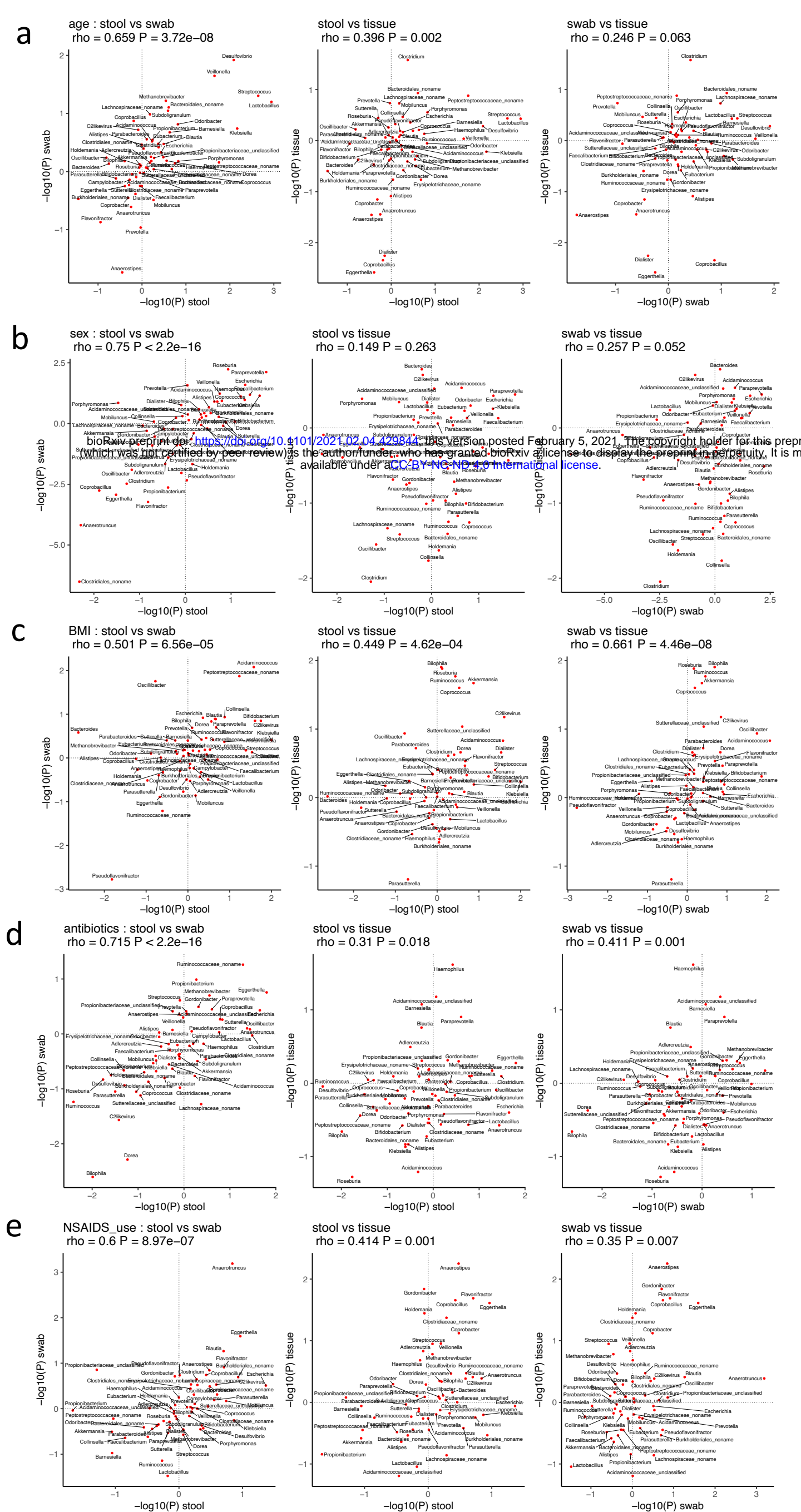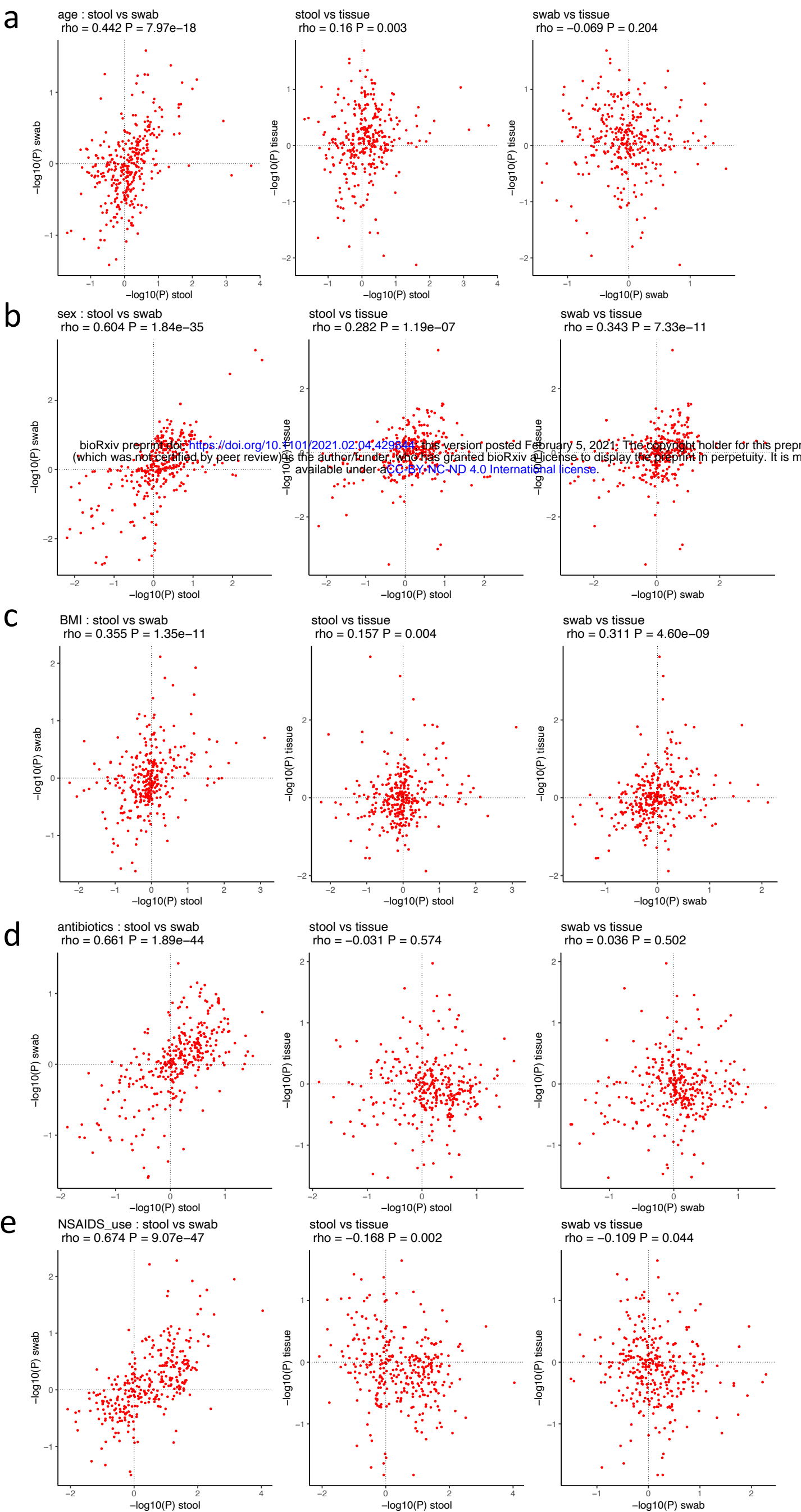Fig 6. Correlations between the functional pathways inference for age (a), sex (b), BMI (c), antibiotics use (d) , NSAIDs use (e) and between pairwise sample types. The axes were the −log10 transformation of p-values from the model2 described in methods.