

A metagenomic DNA sequencing assay that is robust against environmental DNA contamination

Omary Mzava^{1*}, Alexandre Pellan Cheng^{1*}, Adrienne Chang^{1*}, Sami Smalling¹, Liz-Audrey Djomnang Kounatse¹, Joan Lenz¹, Randy Longman², Amy Steadman³, Mirella Salvatore⁴, Manikkam Suthanthiran^{5,6}, John R. Lee^{5,6}, Christopher E. Mason⁷, Darshana Dadhanian^{5,6}, Iwijn De Vlaminc^{1†}

Affiliations:

¹Nancy E. and Peter C. Meinig School of Biomedical Engineering, Cornell University, Ithaca, New York, USA

²Jill Roberts Center for IBD, Weill Cornell Medicine, Division of Gastroenterology and Hepatology, New York, NY, USA

³Global Health Labs, Bellevue, WA, USA

⁴Division of Public Health Programs, Department of Medicine, Weill Cornell Medicine, New York, New York, USA

⁵Division of Nephrology and Hypertension, Department of Medicine, Weill Cornell Medicine, New York, NY, 10065, USA

⁶Department of Transplantation Medicine, New York Presbyterian Hospital–Weill Cornell Medical Center, New York, NY, 10065, USA

⁷Department of Physiology and Biophysics, Weill Cornell Medical College, New York City, NY, USA

*These authors contributed equally.

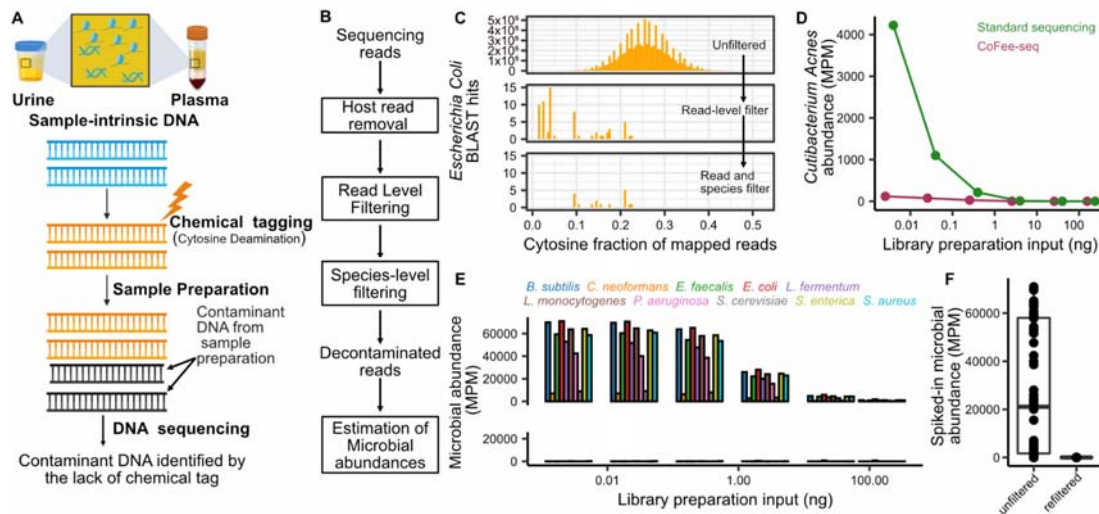
†Corresponding author at: vlaminck@cornell.edu

ABSTRACT (150 words)

Metagenomic DNA sequencing is a powerful tool to characterize microbial communities but is sensitive to environmental DNA contamination, in particular when applied to samples with low microbial biomass. Here, we present contamination-free metagenomic DNA sequencing (Coffee-seq), a metagenomic sequencing assay that is robust against environmental contamination. The core idea of Coffee-seq is to tag the DNA in the sample prior to DNA isolation and library preparation with a label that can be recorded by DNA sequencing. Any contaminating DNA that is introduced in the sample after tagging can then be bioinformatically identified and removed. We applied Coffee-seq to screen for infections from microorganisms with low burden in blood and urine, to identify COVID-19 co-infection, to characterize the urinary microbiome, and to identify microbial DNA signatures of inflammatory bowel disease in blood.

INTRODUCTION

Metagenomic DNA sequencing is a routinely used tool to characterize the genetic makeup and species composition of microbial communities. In addition, metagenomic DNA sequencing of clinical isolates is increasingly used for unbiased detection of microbial infection. Nonetheless, sample contamination by environmental DNA plagues these assays. DNA contamination unavoidably occurs to a degree during the process of sample preparation for DNA sequencing and is particularly problematic for samples that have a low biomass of microbial DNA that can easily be overwhelmed by contaminating DNA¹⁻³.



46
47 **Figure 1. Coffee-seq proof-of-principle.** **A)** Experiment workflow. Tagging of sample-intrinsic DNA by
48 bisulfite DNA treatment is performed directly on urine or plasma. Contaminating DNA introduced after the
49 tagging step is identified based on lack of cytosine conversion. **B)** Bioinformatics workflow. **C)** Representative
50 example of the cytosine fraction of mapped reads in an unfiltered (top) dataset, a read-level filtered dataset
51 (middle) and a fully filtered dataset (bottom). **D)** Number of reads assigned to *Cutibacterium acnes* (common
52 environmental DNA contaminant) in Φ X174 DNA after conventional sequencing (green) and Coffee-seq
53 (purple). **E)** Deliberate contamination assay. Detection of known contaminants before (top) and after (bottom)
54 filtering. **F)** Number of reads assigned to contaminants.

55

56 Multiple solutions have been proposed to overcome the impact of DNA contamination on low
57 biomass metagenomic sequencing. DNA contamination can be avoided to an extent by processing
58 samples in a clean room facility^{4,5}. However, this approach does not avoid contaminant DNA
59 present in reagents. Other approaches are based on batch-correction algorithms that identify
60 microbial species detected in negative controls^{5,6}. These methods however, tend to overcorrect,
61 eliminate sample-intrinsic species that are also common DNA contaminants, and make the incorrect
62 assumption that sample contamination is perfectly reproducible across all samples in a batch. Here,
63 we describe Contamination-Free metagenomic sequencing (Coffee-seq), a metagenomic
64 sequencing method that is robust against DNA contamination. Coffee-seq tags sample-intrinsic,
65 non-contaminant DNA, before DNA isolation with a chemical label that can be recorded via DNA
66 sequencing. Contaminating DNA that is introduced in the sample after this initial tagging step can
67 then be identified and eliminated. Several biochemistries can be envisioned for the initial DNA
68 tagging step. Here, we implement deamination of unmethylated cytosines via bisulfite salt treatment
69 of DNA. This chemistry does not require the use of enzymes or DNA oligos and can be applied
70 directly to clinically relevant samples, such as blood and urine, as demonstrated in this work. We
71 present an analysis of the technical performance of Coffee-seq and describe proof-of-principle
72 applications of Coffee-seq to identify viral and bacterial COVID-19 co-infection from blood, to screen
73 for urinary tract infection (UTI), to characterize the urinary microbiome, to screen for infections with
74 low burden and prevalence in the blood of patients that presented with respiratory symptoms at
75 outpatient clinics in Uganda, and to identify microbial DNA signatures in the blood of patients with
76 inflammatory bowel disease (IBD).

77

78 **Coffee-seq working principle**

79 For the practical implementation of Coffee-seq, we tag DNA by bisulfite salt-induced conversion of
80 unmethylated cytosines to uracils (**Fig. 1A**). Uracils created by bisulfite treatment are converted to
81 thymines in subsequent DNA synthesis steps that are part of the DNA sequencing library
82 preparation. After DNA sequencing, contaminating DNA introduced after tagging can then be

83 identified based on the lack of cytosine conversion. Bisulfite conversion does not require the use of
84 commercial enzymes or oligos that are a frequent source of DNA contamination, and we found that
85 it can be applied directly to the original sample, before DNA isolation. We developed a
86 bioinformatics procedure to differentiate sample-intrinsic microbial DNA, contaminant microbial
87 DNA, and host-specific DNA after Coffee-seq tagging (**Fig. 1B**, Methods). This procedure consists
88 of three steps. First, host cfDNA is removed via mapping and k-mer matching. Given that CpG
89 dinucleotides are heavily methylated in the human genome and rarely in microbial genomes,
90 sequences containing CG dinucleotides are also removed. Second, remaining sequences that
91 consist of more than three cytosines, or one cytosine-guanine dinucleotide are flagged and removed
92 as likely contaminants. Last, a species-level filtering step is performed to remove any remaining
93 reads that primarily originate from C-poor regions in the reference genome (**Fig. 1C**, Methods).

94 We devised two assays to test the principle of Coffee-seq. First, we applied Coffee-seq and
95 conventional DNA sequencing to samples of sheared Φ X174 DNA (New England Biolabs,
96 #N3021S) with variable biomass (0.0025 ng, 0.025 ng, 0.25 ng, 2.5 ng, 26 ng, and 155 ng for
97 Coffee-seq; 0.004 ng, 0.04 ng, 0.4 ng, 4 ng, 35 ng, and 240 ng for standard cfDNA sequencing). We
98 first quantified the abundance of *Cutibacterium acnes* (*C. acnes*), which is a frequent member of the
99 normal skin flora and is routinely identified as a contaminant in DNA sequencing⁷. We observed an
100 increase in *C. acnes* abundance with decreasing input biomass, as expected given that samples
101 with a lower biomass are more susceptible to environmental contamination (**Fig. 1C**). We found that
102 despite a ~30% lower biomass at the beginning of library preparation for the Coffee-seq samples,
103 far fewer *C. acnes* reads were present after Coffee-seq filtering (4223.8 and 119.5 MPM in the
104 highest biomass samples, 1.48 and 0 MPM in the lowest biomass samples, before and after Coffee-
105 seq filtering respectively; **Fig. 1D**).

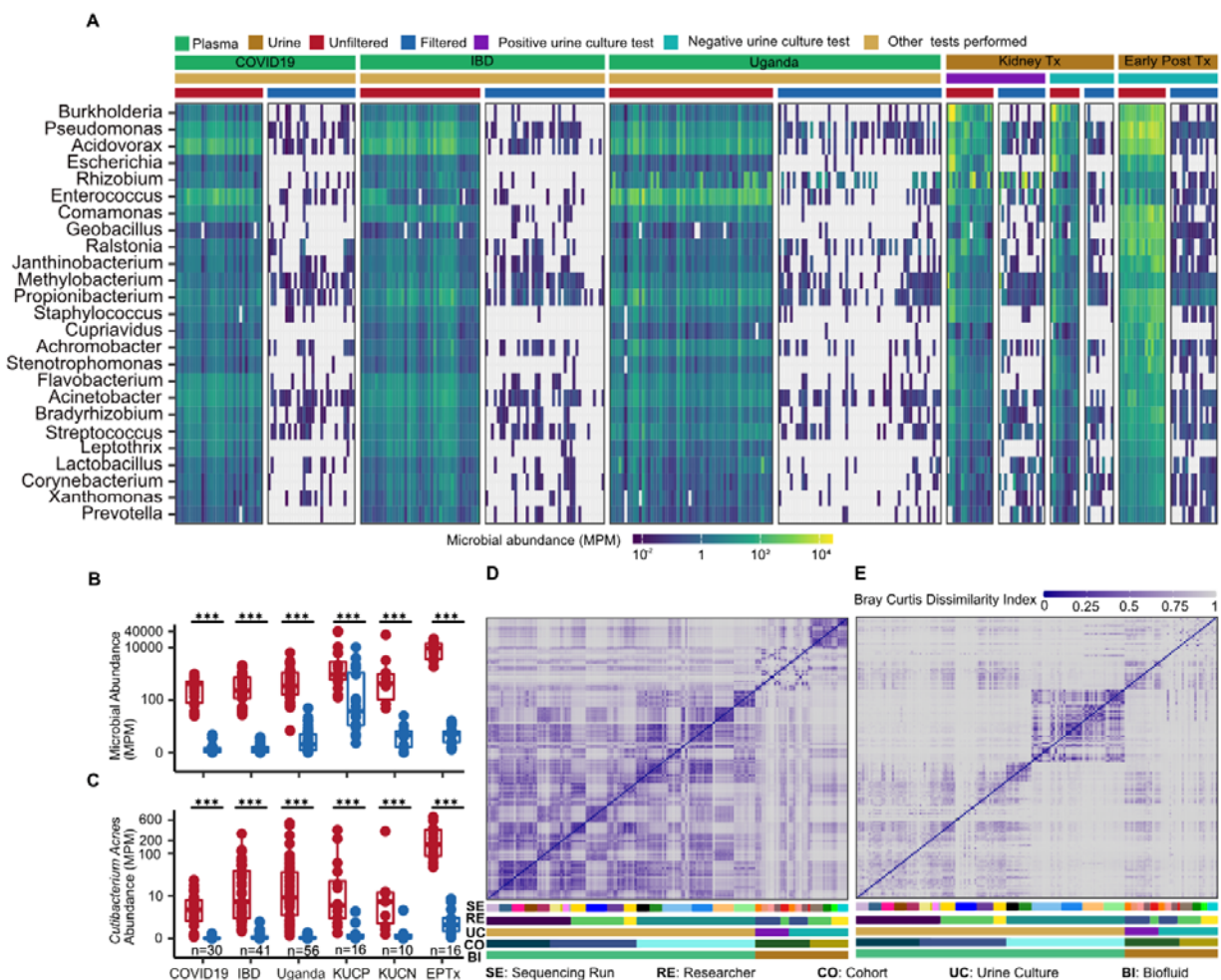
106 Second, we performed Coffee-seq on sheared Φ X174 DNA samples with variable biomass (0.0025-
107 155 ng; **Fig. 1E**) which we spiked after Coffee-seq tagging with 1 ng of sheared DNA from a well-
108 characterized community of microbes to simulate microbial DNA contamination (10 species; Zymo
109 Research, #D6305). Before applying the Coffee-seq bioinformatics filter, we observed a negative
110 correlation between the Φ X174 DNA input biomass and the relative number of reads from the spike-
111 in community, as expected (Pearson's $R = -0.54$, p -value = 6.5×10^{-6} ; Spearman's $\rho = -0.82$, p -value
112 = 6.3×10^{-16} ; **Fig. 1E**). After applying the Coffee-seq filter, we observed an average percent decrease
113 of 99.8% of molecules mapping to species of the spike-in community (**Fig. 1F**). Sequences mapping
114 to *Escherichia coli* (*E. coli*) were the most abundant after filtering (58.89%). Given that Φ X174
115 genomic DNA is isolated after phage propagation in *E. coli* culture, we reasoned that these
116 remaining reads were likely intrinsic to the original sample. Together, these experiments
117 demonstrate the effectiveness of Coffee-seq for the detection and removal of DNA contaminants.

118 **Application of Coffee-seq to cell-free DNA in blood and urine**

119 Cell-free DNA (cfDNA) in blood and urine has emerged as a useful analyte for the diagnosis of
120 infection⁸⁻¹⁵. Metagenomic cfDNA sequencing can identify a broad range of potential pathogens with
121 high sensitivity. Yet, because of the low biomass of microbial-derived cfDNA in blood and urine,
122 metagenomic cfDNA sequencing is highly susceptible to environmental contamination, limiting the
123 specificity of metagenomic cfDNA sequencing for pathogen identification.

124 To assess the performance of Coffee-seq in metagenomic cfDNA sequencing, we assayed a total of
125 169 cfDNA isolates (42 urine, 127 plasma) collected from five groups of subjects: **1**) 26 urine
126 samples from a cohort of kidney transplant patients with and without UTI (16 UTI positive, 10 UTI
127 negative; "kidney transplant cohort"), **2**) 16 urine samples collected early after transplantation from
128 10 kidney transplant patients that received a ureteral stent at the time of transplantation (samples

129 were collected pre-stent and post-stent removal for 5 of the 10 patients; “early post-transplant
 130 cohort”), **3)** 56 plasma samples from a cohort of 44 patients presenting with respiratory symptoms at
 131 outpatient clinics in Uganda (28 sputum positive for Tuberculosis [TB], 16 sputum negative for TB;
 132 “Uganda cohort”), **4)** 41 plasma samples from a cohort of 32 patients diagnosed with IBD (16
 133 patients with Crohn’s disease, 16 patients with ulcerative colitis; “IBD cohort”), and, **5)** 30 plasma
 134 samples from a cohort of 14 patients hospitalized with COVID-19 (“COVID-19 cohort”; see **Table S1**
 135 and Supplementary Information for details on the patients and samples included).



136

137 **Figure 2. Coffee-Seq applied to cell-free DNA in urine and plasma. A)** Microbial abundance of 25 most
 138 abundant common contaminant genera (selected from the 68 genera⁴) before and after Coffee-seq filtering in
 139 plasma and urine from five independent subject cohorts (Tx = transplant). Total abundance of all contaminant
 140 genera **B)** and *C. acnes* **C)** before and after Coffee-seq filtering (KUCP = Kidney Transplant cohort with
 141 positive urine culture, KUCN = Kidney Transplant cohort with negative urine culture, EPTx = Early Post
 142 Transplant cohort). Bray-Curtis dissimilarity index before **D)** and after **E)** filtering. Samples are organized by:
 143 sequencing batch, researcher performing the experiment, cohort, and biofluid. *** p-value < 0.001

144 We performed Coffee-seq for all samples and obtained an average of 46.5 ± 23.6 million paired-end
 145 reads per sample. We detected and quantified the abundance of 68 genera that have been reported
 146 as frequent DNA contaminants in multiple independent studies (summarized in Ref. 4; **Fig 2A**, 49
 147 of these genera detected in at least one sample). We found that 76% of these genera were completely
 148 removed from all samples after Coffee-seq filtering. We calculated the total number of molecules
 149 from all contaminant genera and observed an up to 3 orders of magnitude reduction after Coffee-
 150 seq filtering (reduced by a factor of 7.5, 1711.2, 177.6, 548.3, 547.2 for the kidney transplant cohort,
 151 early post-transplant cohort, Uganda cohort, IBD cohort, and COVID-19 cohort, respectively; **Fig.**

152 **2B**). We investigated the impact of Coffee-seq filtering on removing reads originating from the skin
153 contaminant *C. acnes* (**Fig. 2C**). *C. acnes* was detected in all samples and completely removed
154 from 50 samples by Coffee-seq filtering. In the remaining samples, we observed an up to 2 orders of
155 magnitude reduction of *C. acnes* reads.

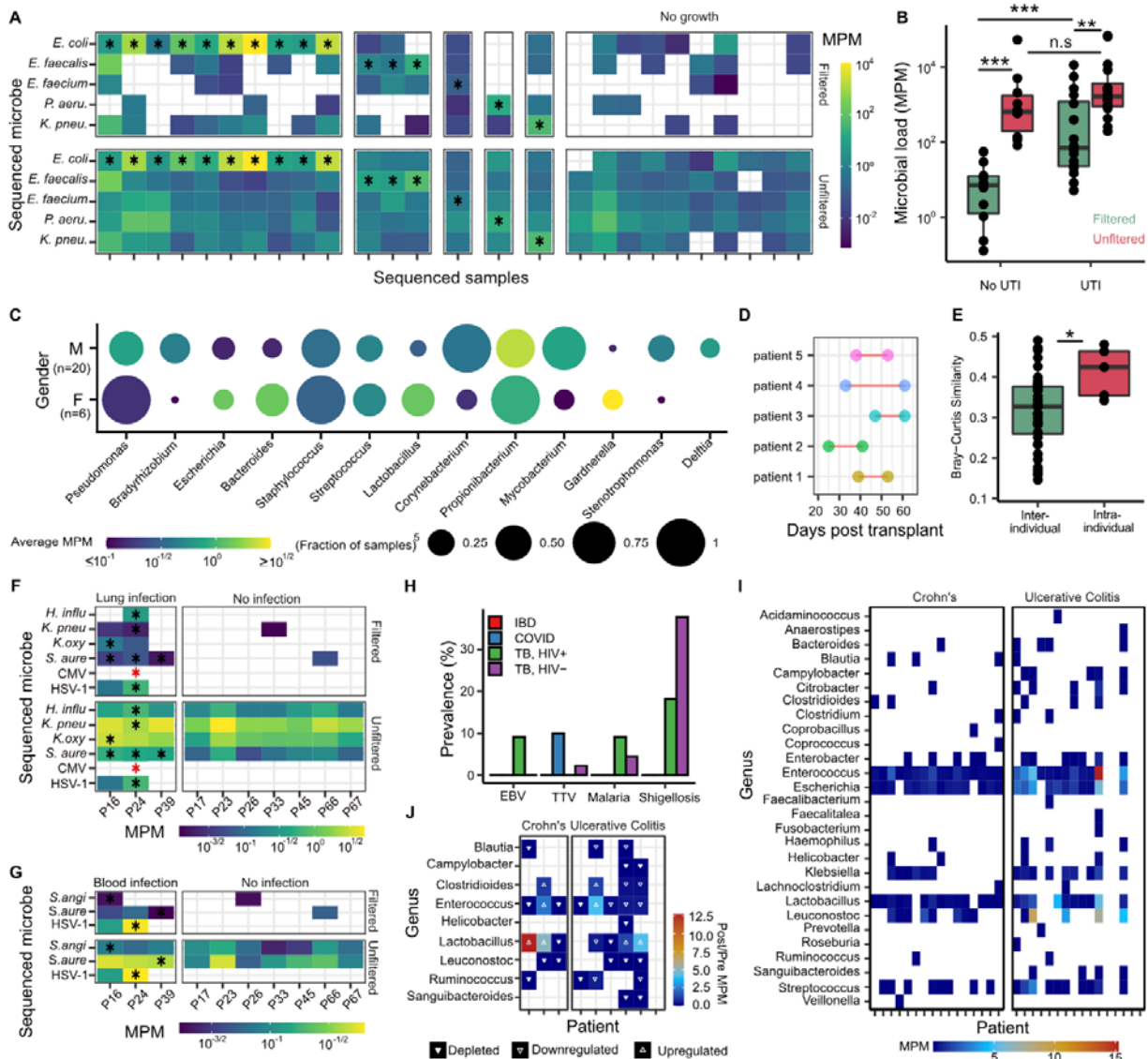
156 We next evaluated the utility of Coffee-seq to correct for batch effects and to reveal true differences
157 in microbiome profiles for different patient groups. To this end, we calculated the Bray-Curtis
158 Dissimilarity Index for all clinical samples included in this study and sorted the datasets based on
159 the following parameters: **1**) sequencing run, **2**) operator, **3**) urine culture test, **4**) study cohort, and
160 **5**) biofluid type. Before Coffee-seq filtering, we observed a high similarity for samples assayed in the
161 same experimental batches (**Fig. 2D**). Coffee-seq filtering removed these batch effects and
162 revealed distinct cohort-specific microbiome profiles. Most notably, we observed distinct plasma
163 microbiome profiles for plasma samples from the Uganda cohort (**Fig. 2E**). These results
164 demonstrate that Coffee-seq directly applied to biofluids leads to a dramatic decrease in
165 experimental noise and bias due to DNA contamination.

166 **Coffee-seq enables to screen for UTI and to characterize the urine microbiome**

167 The healthy urinary tract was long believed to be sterile^{16,17}, but this picture was challenged with
168 recent advances in urine culture techniques that have identified bacteria in the urinary tract of both
169 males and females¹⁸. Yet many microbes are difficult to cultivate *in vitro*, and bacterial culture can
170 also be sensitive to contamination¹⁹. Therefore, comprehensive and accurate characterization of
171 species colonizing the urinary microbiome is still lacking.

172
173 We reasoned that Coffee-seq could provide insight into the composition of the urine microbiome
174 with both high sensitivity and specificity. We first applied Coffee-seq to 26 urine samples from 23
175 kidney transplant patients with and without infection of the urinary tract as determined by
176 conventional urine culture (16 UTI positive [*Enterococcus faecalis*: n=3; *Enterococcus faecium*: n=1;
177 *Escherichia coli*: n=10; *Klebsiella pneumoniae*: n=1; *Pseudomonas aeruginosa*: n=1] and 10 UTI
178 negative). Coffee-seq consistently identified microbial cfDNA from species reported by urine culture
179 (16/16 UTI positive samples; **Fig. 3A**). Coffee-seq also identified two *Corynebacterium* species
180 (*Corynebacterium jeikeium* and *Corynebacterium urelyticum*) in one sample from a UTI positive
181 patient (*E.coli*) with culture confirmed *Corynebacterium* co-infection. In addition, we found that
182 samples from UTI positive patients had a significantly higher burden of total microbial DNA
183 compared to samples from UTI negative patients (1451.8 ± 3024.7 MPM and 12.8 ± 17.6 MPM,
184 respectively in the filtered samples; p-value = 1.1×10^{-5} , Wilcoxon test; **Fig. 3B**). Conventional
185 metagenomic sequencing (without Coffee-seq filtering) detected uropathogens with equal sensitivity
186 but suffered from poor specificity: DNA from common uropathogens not identified by culture was
187 detected in many samples, albeit with low abundance, including in samples from patients without
188 UTI. We conclude that the improved specificity of Coffee-seq allows for more accurate
189 characterization of co-infection networks in the scope of UTIs, and more accurate characterization
190 of the normal urine microbiome in the absence of UTIs. It is important to note that two common skin
191 microbes, *C. acnes* and *Staphylococcus epidermis*, were found in most samples (23/26 samples).
192 While these two species have been shown to cause UTIs^{20,21}, they may also have been introduced
193 as contaminants at the time of urine collection, which underscores an important limitation of Coffee-
194 seq: Coffee-seq is not robust against contamination that occurs before the tagging step.

195



196
 197 **Figure 3. Application of Coffee-seq to plasma and urine.** **A)** Heatmap of abundance of species (molecules
 198 per million, MPM) identified in patients with and without UTI, before and after application of Coffee-seq filter.
 199 **B)** Boxplot of the relative number of microbe-derived molecules (MPM) in samples from patients with and
 200 without UTI, before and after Coffee-seq filtering. **C)** Dot plot of the most abundant genera in urine from male
 201 and female kidney transplant recipients. **D-E)** Boxplot showing Bray-Curtis similarity index (as defined in **D)**
 202 of the urine microbiome within individual patients and between patients before and after stent removal. **F-G)**
 203 Heatmaps of the abundance of species identified in plasma from COVID-19 patients with and without culture
 204 confirmed **F)** lung and **G)** blood infection, before and after application of Coffee-seq filter (red * indicates
 205 detection by sputum culture only). Red boxes indicate positive culture tests. **H)** Barplot of the prevalence of
 206 Epstein-Barr Virus (EBV), Torque teno virus (TV), Malaria, or Shigellosis pathogens in different patient
 207 cohorts. **I)** Heatmap of the abundance of species identified in matched stool and plasma cfDNA samples in
 208 patients diagnosed with Crohn's disease or ulcerative colitis. **J)** Heatmap of the change in abundance of gut
 209 specific bacteria before and after treatment. (Black * in panels A, F, and G indicates agreement with urine,
 210 respiratory and blood culture, respectively).
 211

212 To explore the effect of gender on the urine microbiome, we analyzed isolates from culture
 213 confirmed UTI negative patients (n=26) from the kidney transplant (n=10) and early post-transplant
 214 (n=16) study cohorts (5 female, 14 male). This analysis yielded a small, but statistically insignificant,
 215 difference in total microbial load for male versus female patients (**Fig. S1**). We also observed that a

216 subset of the most abundant genera was found in both male and female samples, with a marked
217 variation in number of samples and abundances (**Fig. 3C**).

218 Studies investigating the temporal dynamics of urine microbiome in individuals can benefit from the
219 high sensitivity and specificity achieved with our assay. We applied Coffee-seq to paired urine
220 samples obtained from 5 kidney transplant patients collected at two time points before and after
221 ureteral stent removal (**Fig. 3D**). We compared the similarity of microbial composition between
222 samples from the same patient (intra-individual) and between different patients (inter-individual) at
223 different sampling points and observed that the microbial composition remained more similar in the
224 same patient (**Fig. 3E**) than between different patients, supporting the utility of Coffee-seq to
225 measure subtle dynamics in urine microbiome composition (Mean Bray-Curtis Similarity: 0.41 ± 0.06
226 and 0.317 ± 0.09 respectively, p -value = 3.1×10^{-2} , Wilcoxon test).

227 **Coffee-seq identifies bacterial and viral co-infection of COVID-19 from blood**

228 The COVID-19 pandemic is an unprecedented human health crisis. Viral or bacterial co-infection
229 occurs in roughly 4% of hospitalized COVID-19 patients but can occur in up to 30% of COVID-19
230 patients admitted to the intensive care unit²². Co-infection has been associated with longer fever
231 duration, and increased admittance to the intensive care unit and ventilation treatment²³. We
232 reasoned that Coffee-seq may offer sensitive detection of bacterial co-infection in COVID-19
233 patients with improved specificity over conventional metagenomic sequencing assays.

234 We applied Coffee-seq to 30 plasma samples from 14 patients with COVID-19 collected as part of a
235 clinical study aimed at identifying predictors of disease severity. Respiratory and blood cultures
236 were obtained as part of standard clinical care. Three patients (P16, P24, P39) tested positive for
237 blood borne infection and respiratory tract infection, while all other patients were not diagnosed with
238 COVID-19 co-infection. Coffee-seq identified the causative pathogen in 3/3 blood infection cases
239 and 7/8 respiratory infection cases (**Fig. 3F-G**). Conventional metagenomic sequencing (without
240 Coffee-seq filtering) was equally sensitive to these pathogens but was limited by specificity (**Fig. 3F-**
241 **G**). Of interest, while we did not obtain plasma collected the day of infection for P24, we identified
242 cfDNA originating from *K. pneumoniae* and *Haemophilus influenzae*, for which the patient tested
243 positive four days later. While further investigation is necessary to resolve discrepancies between
244 positive culture results and microbial cfDNA detection, these results suggest that Coffee-seq may be
245 able to identify cases of infection earlier than traditional culture methods, and with improved
246 specificity compared to conventional metagenomic sequencing techniques.

247 **Coffee-seq identifies bacterial and viral infections with low prevalence and low microbial** 248 **burden**

249 Neglected tropical diseases significantly impact the public health and economies of low-income
250 countries. Treatments exist for many of these diseases, but development and deployment of reliable
251 diagnostic tests has been slow²⁴. We reasoned that Coffee-seq could be used to screen for
252 infections with low prevalence and low microbial burden.

253 We applied Coffee-seq to 56 plasma samples from 44 individuals who presented with symptoms of
254 respiratory illness at outpatient clinics in Uganda (28 sputum positive tuberculosis, 16 sputum
255 negative tuberculosis). Nine of these individuals were HIV positive at the time of sample collection.
256 We mined the data to determine the prevalence of infections endemic to Uganda and compared
257 with results obtained for plasma samples collected from subjects that live in North America (54
258 plasma samples from the IBD cohort; 30 plasma samples from the COVID-19 cohort). We screened
259 the samples for Epstein-Barr virus, Torque Teno virus, and pathogens associated with malaria

260 (*Plasmodium vivax* and *P. falciparum*), and shigellosis (*Shigella sonnei*, *S. dysenteriae*, *S. boydii*,
261 and *S. flexneri*). These pathogens were found at varying rates in samples from the Uganda cohort
262 (**Fig. 3H**): malaria (3/44), Epstein-Barr virus (1/44), shigellosis (19/44), and torque teno virus (1/44),
263 but not in the IBD cohort. Torque teno virus, which has previously been reported to be elevated in
264 immunocompromised patients⁸, was identified in 3/30 COVID-19 patient samples, all from patients
265 who had received a bone marrow transplant prior to sample acquisition.

266 **Coffee-seq identifies signatures of bacterial translocation from the gastrointestinal tract**

267 Bacterial translocation of intestinal microbes through mucosal membranes is believed to be a
268 normal phenomenon, but has been found to occur more frequently in patients experiencing gut flora
269 disruption^{25,26}. In patients with inflammatory bowel disease, gut vascular barrier disruption has been
270 linked to increased intestinal permeability and subsequent microbial translocation across the
271 mucosal membrane^{27,28}. The translocation of gut bacteria and their products to extraintestinal sites
272 can result in systemic inflammation, resulting in autoimmune or other non-infectious diseases.
273 Detecting signatures of translocation is therefore important but difficult in view of the low abundance
274 of microbial DNA due to translocation in blood.

275 To identify signatures of bacterial translocation, we compared whole genome shotgun sequencing of
276 fecal samples from 32 patients (Crohn's n=16, ulcerative colitis, n=16) to matched plasma cfDNA
277 samples assayed using Coffee-seq. We first quantified bacterial species identified in matched fecal
278 and plasma samples (**Fig. 3I**). We identified cfDNA derived from gut-specific microbes in all patient
279 samples, though to a much greater extent in individuals with ulcerative colitis (1.40±1.4 vs
280 6.82±10.6 MPM of gut specific bacteria for Crohn's disease and ulcerative colitis, respectively). To
281 investigate the effects of treatment on bacterial translocation, we collected additional stool and
282 plasma samples from nine patients (Crohn's n=3, ulcerative colitis n=6) after treatment initiation and
283 performed whole genome shotgun sequencing of stool and Coffee-seq on plasma cfDNA. We
284 quantified the relative abundance of gut-specific bacterial species before and after treatment and
285 found that the burden of cfDNA decreased for most bacterial species (28/36) following treatment,
286 which may be explained by a reduction in the degree of bacterial translocation with treatment (**Fig.**
287 **3J**). Of interest, out of seven subjects for which we detected *Lactobacillus* before treatment, five
288 displayed an increase in *Lactobacillus* species burden in blood after treatment (up to 12.7-fold
289 increase after treatment and an average of 3.36-fold MPM increase after treatment across all
290 samples). *Lactobacillus* has been shown to promote gastrointestinal barrier function, protecting the
291 gut from pathogenic bacteria and preventing inflammation²⁸. For bacterial species besides
292 *Lactobacillus*, we find an average of 0.3-fold MPM reduction after treatment. These preliminary
293 results support the use of Coffee-seq to identify subtle signatures of bacterial translocation in the
294 blood.

295 **DISCUSSION**

296 We report Coffee-seq, a method for metagenomic DNA sequencing that is robust against DNA
297 contamination. In contrast to prior methods for the management of DNA contamination that have
298 relied on algorithmic batch correction or the use of known-template or no-template controls, Coffee-
299 Seq uses a physical labeling technique to differentiate sample-intrinsic DNA from contaminating
300 DNA. The principle of Coffee-seq has the potential for broad application in contexts where
301 metagenomic analyses of isolates with low biomass of microbial DNA are required. In this proof-of-
302 principle study, we have explored applications of Coffee-seq to quantify microbial cell-free DNA in
303 human biofluids. Metagenomic sequencing of microbial cell-free DNA in blood or urine is a highly
304 sensitive approach to screen for a broad range of viral or bacterial pathogens, but because of the
305 low biomass of microbial DNA in blood and urine this method is highly susceptible to DNA

306 contamination leading to a high false positive rate. We implemented Coffee-seq tagging of cell-free
307 DNA in plasma and urine by bisulfite-induced deamination of unmethylated cytosines and show that
308 this approach reduces background signals from common contaminants by up to three orders of
309 magnitude. Coffee-seq thereby dramatically improves the specificity of metagenomic cfDNA
310 analyses, opening up a broad range of applications, e.g. infectious disease with low microbial
311 burden or syndromes that are accompanied by subtle changes in the plasma or urine microbiome.

312 In its current implementation, Coffee-seq has several limitations. First, Coffee-seq is only robust
313 against DNA contamination introduced after the labeling step. We implemented Coffee-seq tagging
314 directly on biofluids, which allowed us to identify contaminants introduced during DNA isolation or
315 library preparation but not during the sample collection or isolation of the plasma from whole blood.
316 Second, the specific labeling strategy we have implemented here inherently modifies the DNA
317 sequence and thereby limits the resolution of sequence-based analyses. Alternative
318 implementations of contamination-free sequencing that do not introduce sequence alterations can
319 be considered. Last, the principles introduced here can be adopted for molecular assays beyond
320 whole genome sequencing, including amplicon sequencing assays, e.g. 16S rRNA profiling, or PCR
321 assays.

322 **METHODS**

323 **Study Cohort and sample collection:**

324 *Uganda cohort and sample collection*

325 Forty-four plasma samples were collected from individuals seeking tuberculosis treatment in
326 Uganda. Briefly, peripheral blood was collected in Streck Cell-Free BCT (Streck #230257) and
327 centrifuged at 1600 x g for 10 minutes. Plasma was stored in 1 mL aliquots at -80°C. The study was
328 approved by the Makerere School of Medicine Research and Ethics Committee (protocol 2017-020).
329 All patients provided written informed consent.

330 *IBD cohort sample collection*

331 Peripheral blood samples were collected under IRB approved protocol (1806019340) at the Jill
332 Roberts Center for IBD at Weill Cornell Medicine. PBMCs and plasma were fractionated using a
333 Ficoll-Hypaque gradient.

334 *Stool sample collection*

335 DNA from fecal samples was isolated using the MagAttract PowerMicrobiome DNA/RNA kit with
336 glass beads (Qiagen, Germany). Metagenomic libraries were prepared using the NEBNext Ultra II
337 for DNA Library Prep kit (New England Biolabs, Ipswich, MA) following the manufacturer's protocol.
338 The DNA library was sequenced on an Illumina HiSeq instrument using a 2x150 paired-end
339 configuration in a high output run mode.

340 *COVID-19 cohort sample collection*

341 Samples were collected as part of an observational study among individuals with COVID-19^{29,30} that
342 were treated at New York Presbyterian Hospital and Lower Manhattan Hospitals, Weill Cornell
343 Medicine. The study was approved by the Institutional Review Board of Weill Cornell Medicine (IRB
344 20-03021645), and informed consent was obtained from all participants.

345 *UTI cohort sample collection*

346 Twenty six urine samples were collected from 23 kidney transplant recipients who received care at
347 New York Presbyterian Hospital–Weill Cornell Medical Center. The study was approved by the Weill
348 Cornell Medicine Institutional Review Board (protocols 1207012730). All patients provided written
349 informed consent. Patients provided urine specimens using a clean-catch midstream collection
350 protocol. The urine specimen was centrifuged at 3000 x g for 30 minutes and supernatant was
351 stored as 1 mL of 4 mL aliquots.

352 *Early post transplant sample collection*

353 Urine specimens collected within 10 ± 5 days of ureteral stent removal from patients who agreed to
354 participate in the WCM IRB approved protocol # 20-01021269 were included in this study. Urine
355 specimens were collected within 47 ± 11 days post-kidney transplantation. The presence of UTI was
356 excluded by a negative urine culture and the absence of pyuria. This study was approved by the
357 Weill Cornell Medicine Institutional Review Board (protocol 20-01021269).

358 *Definition of Positive and Negative urine culture for the UTI and Early post-transplant cohorts*

359 A positive urine culture was defined as a culture growing an organism identified to at least the genus
360 level ($\geq 10,000$ cfu/mL). A urine culture was defined as negative when either no organism was
361 isolated in culture (< 1000 cfu/mL) or the organism was unidentified to either the genus or species
362 level (i.e., unidentified) and the colony count was $< 10,000$ cfu/mL.

363 **Coffee-seq in plasma.** An aliquot of 520 μ L of plasma was centrifuged at 14,000 RPM for 10
364 minutes at 12°C to pellet cellular debris. The supernatant was transferred to a new 1.5 mL tube and
365 the final volume was brought up to 1000 μ L with PBS. The solution was heated to 98°C for 10
366 minutes and mixed at 1000 RPM to coagulate the albumin present in plasma. The solution was then
367 centrifuged at 4000 RPM for 10 minutes. 500 μ L of supernatant was transferred to 15 mL falcon
368 tube containing 3.25 mL of ammonium bisulfite solution (Zymo Research, product #5030) and
369 shaken in a thermomixer at 98°C for 10 minutes (15s on/30s off). Samples were then transferred to
370 a thermomixer at 54°C for 60 minutes (15s on/30s off). Then, cfDNA extraction was performed using
371 the QIAamp Circulating Nucleic Acid Kit using the 4-mL plasma protocol (Qiagen, product #55114).
372 Prior to DNA elution, 200 μ L of L-Desulphonation buffer (Zymo Research, product #5030) was
373 added to the columns for 15 minutes, followed by two washes with 200 μ L absolute ethanol. DNA
374 was then eluted according to manufacturer recommendations, and single-stranded library
375 preparation is performed (Claret Biosciences, product #CBS-K150B). Libraries were then
376 sequenced on an Illumina sequencer.

377 **Coffee-seq in urine.** An aliquot of 520 μ L of urine was centrifuged at 14,000 RPM for 5 minutes to
378 pellet cellular debris. 500 μ L of supernatant was transferred to a new 15 mL falcon tube containing
379 3.25 mL of ammonium bisulfite solution (Zymo Research, product #5030) and heated to 98°C for 10
380 minutes. Samples were then kept at 54°C for 60 minutes. Then, cfDNA extraction was performed
381 using a commercially available column-based kit (Norgen Biotek, product #56700). Prior to DNA
382 elution, 200 μ L of L-Desulphonation buffer (Zymo Research, product #5030) was added to the
383 columns for 20 minutes, followed by two washes with 200 μ L absolute ethanol. DNA was then
384 eluted according to manufacturer recommendations, and single-stranded library preparation was
385 performed (Claret Biosciences, product #CBS-K150B). Libraries were then sequenced on an
386 Illumina sequencer.

387 **Alignment to the human genome.** Adapter and low quality bases from the reads were trimmed
388 using BBDuk³¹ and aligned to the C-to-T and G-to-A converted human genome using Bismark³²
389 (Bismark-0.22.1). PCR duplicates were removed using Bismark.

390 **Depth of coverage.** The depth of sequencing was measured by summing the depth of coverage for
391 each mapped base pair on the human genome after duplicate removal, and dividing by the total
392 length of the human genome (hg19, without unknown bases).

393 **Bisulfite conversion efficiency.** We estimated bisulfite conversion efficiency by quantifying the
394 rate of C[A/T/C] methylation in human-aligned reads (using MethPipe³³ V3.4.3), which are rarely
395 methylated in mammalian genomes.

396 **Metagenomic abundance estimation from sequencing data.** Metagenomic analysis is performed
397 as previously described¹². Specific to Coffee-seq, read-level filtering of contaminants is performed
398 by removing sequenced reads with 4 or more cytosines present, or one methylated CpG
399 dinucleotide (the latter represents unmapped, human-derived molecules). Species-level filtering
400 based on the distribution of mapped reads is carried out by first aligning filtered and unfiltered
401 datasets independently. Cytosine-densities of mapping-coordinates in both datasets are measured
402 using custom scripts, and their distributions are compared using a Kolmogorov-Smirnov test.
403 Significantly different filtered-unfiltered distributions are further processed (D-statistic > 0.1 and p-
404 value < 0.01). Briefly, filtered datasets whose distribution of cytosines at mapped locations is
405 significantly lower than unfiltered datasets have one read removed, and are re-tested for differences
406 in their distribution. If the distributions are more similar (as measured through the same criteria), it is
407 filtered out. This process is repeated until distributions are no longer significantly different, or if all
408 reads are removed. Metagenomic abundances of filtered datasets are estimated using GRAMMy as
409 previously described in Ref 12. Microbial abundance in downstream analyses was quantified as
410 Molecules Per Million reads (MPM).

$$411 \quad MPM = \frac{\text{Adjusted Blast hits} \times 10^6}{\text{Total Trimmed Reads}}$$

412 **Identification of translocated gut bacteria in plasma**

413 Fecal shotgun metagenomic data for 41 samples was obtained from 32 patients diagnosed with
414 inflammatory bowel disease (IBD). Low-quality bases and Nextera-specific sequences were
415 trimmed (Trim Galore). Reads were aligned (Bowtie2³⁴) against the human references (UCSC
416 hg19). Unaligned reads were extracted and assembled with metaSPAdes³⁵ and classified with
417 Kaiju³⁶.

418 Paired cfDNA samples were filtered as previously described and aligned to the assembled reads
419 with Bismark. Mapped reads with a minimum quality score of 15 were extracted and filtered for gut-
420 specific microorganisms identified by The Human Gut Microbiome Atlas³⁷.

421 **Statistical analysis**

422 All statistical methods were performed in R version 4.0.5. Groups were compared using a two-sided
423 Wilcoxon Rank Sum test. Boxes in the boxplots indicates 25th and 75th percentile, the band in the
424 box indicated the median and whiskers extend to 1.5 x Interquartile Range (IQR) of the hinge.

425 **Code and Data Availability**

426 All scripts used in this study are available at <https://github.com/omrmzv/CoffeeSeq>. Φ X174 DNA
427 sequencing data used in the proof of principle experiments has been deposited in NCBI's Sequence
428 Read Archive (SRA) under Bioproject ID (PRJNA782310). Sequencing data from human plasma
429 cfDNA will be deposited in the database of Genotypes and Phenotypes (dbGaP)

430 **ACKNOWLEDGMENTS**

431 We thank the Cornell Genomics Center for help with sequencing assays, the Cornell Bioinformatics
432 facility for computational assistance and Michael Satlin and Lars Westblade for helpful discussions.
433 A special thanks to Dr. Alfred Andama for his research supervision with the Infectious Disease
434 Research Collaboration (IDRC) for samples collected and characterized in Kampala, Uganda. This
435 work was supported by R01AI146165 (to I.D.V.), R21AI133331 (to I.D.V.), R21AI124237 (to I.D.V.),
436 DP2AI138242 (to I.D.V.), R01AI151059 (to I.D.V., J.R.L., M.S., D.D.), R37 AI051652 (to M.S), a
437 Synergy award from the Rainin Foundation (to I.D.V. and R.L.), a grant from the Bill and Melinda
438 Gates Foundation INV-003145 (to I.D.V.). A.C. was supported by the National Institutes of Health
439 under the Ruth L. Kirschstein National Research Service Award (6T32GM008267) from the National
440 Institute of General Medical Sciences. A.P.C. is supported by a National Sciences and Engineering
441 Research Council of Canada PGS-D3 fellowship.

442

443 **CONFLICTS**

444 IDV, OM, APC and AC have submitted a patent related to the present work. APC, IDV, DD, and JRL
445 are inventors on the patent US-2020-0048713-A1 titled “Methods of Detecting Cell-Free DNA in
446 Biological Samples.” I.D.V. is a member of the Scientific Advisory Board of Karius Inc., Kanvas
447 Biosciences and GenDX. JRL received research support under an investigator-initiated research
448 grant from BioFire Diagnostics, LLC.

REFERENCES

1. Glassing, A., Dowd, S. E., Galandiuk, S., Davis, B. & Chiodini, R. J. Inherent bacterial DNA contamination of extraction and sequencing reagents may affect interpretation of microbiota in low bacterial biomass samples. *Gut Pathog.* **8**, 24 (2016).
2. Weyrich, L. S. *et al.* Laboratory contamination over time during low-biomass sample analysis. *Mol. Ecol. Resour.* **19**, 982–996 (2019).
3. Salter, S. J. *et al.* Reagent and laboratory contamination can critically impact sequence-based microbiome analyses. *BMC Biol.* **12**, 87 (2014).
4. Eisenhofer, R. *et al.* Contamination in Low Microbial Biomass Microbiome Studies: Issues and Recommendations. *Trends Microbiol.* **27**, 105–117 (2019).
5. Davis, N. M., Proctor, D. M., Holmes, S. P., Relman, D. A. & Callahan, B. J. Simple statistical identification and removal of contaminant sequences in marker-gene and metagenomics data. *Microbiome* **6**, 226 (2018).
6. Burnham, P. *et al.* Separating the signal from the noise in metagenomic cell-free DNA sequencing. *Microbiome* **8**, 18 (2020).
7. Danko, D. *et al.* A global metagenomic map of urban microbiomes and antimicrobial resistance. *Cell* **184**, 3376-3393.e17 (2021).
8. De Vlaminck, I. *et al.* Temporal Response of the Human Virome to Immunosuppression and Antiviral Therapy. *Cell* **155**, 1178–1187 (2013).
9. De Vlaminck, I. *et al.* Circulating Cell-Free DNA Enables Noninvasive Diagnosis of Heart Transplant Rejection. *Sci. Transl. Med.* **6**, 241ra77-241ra77 (2014).
10. De Vlaminck, I. *et al.* Noninvasive monitoring of infection and rejection after lung transplantation. *Proc. Natl. Acad. Sci.* **112**, 13336–13341 (2015).
11. Burnham, P. *et al.* Urinary cell-free DNA is a versatile analyte for monitoring infections of the urinary tract. *Nat. Commun.* **9**, 2412 (2018).
12. Cheng, A. P. *et al.* A cell-free DNA metagenomic sequencing assay that integrates the host injury response to infection. *Proc. Natl. Acad. Sci.* **116**, 18738–18744 (2019).
13. Blauwkamp, T. A. *et al.* Analytical and clinical validation of a microbial cell-free DNA sequencing test for infectious disease. *Nat. Microbiol.* **4**, 663–674 (2019).

14. Cheng, A. P. *et al.* Cell-free DNA tissues of origin by methylation profiling reveals significant cell, tissue, and organ-specific injury related to COVID-19 severity. *Med N. Y. N* **2**, 411-422.e5 (2021).
15. Chang, A. *et al.* Measurement Biases Distort Cell-Free DNA Fragmentation Profiles and Define the Sensitivity of Metagenomic Cell-Free DNA Sequencing Assays. *Clin. Chem.* (2021) doi:10.1093/clinchem/hvab142.
16. Wolfe, A. J. & Brubaker, L. “Sterile Urine” and the Presence of Bacteria. *Eur. Urol.* **68**, 173–174 (2015).
17. Hilt, E. E. *et al.* Urine Is Not Sterile: Use of Enhanced Urine Culture Techniques To Detect Resident Bacterial Flora in the Adult Female Bladder. *J. Clin. Microbiol.* **52**, 871–876 (2014).
18. Gottschick, C. *et al.* The urinary microbiota of men and women and its changes in women during bacterial vaginosis and antibiotic treatment. *Microbiome* **5**, 99 (2017).
19. Wade, W. Unculturable bacteria—the uncharacterized organisms that cause oral infections. *J. R. Soc. Med.* **95**, 81–83 (2002).
20. Lowy, F. D. & Hammer, S. M. Staphylococcus epidermidis Infections. *Ann. Intern. Med.* **99**, 834–839 (1983).
21. Boisrenoult, P. Cutibacterium acnes prosthetic joint infection: Diagnosis and treatment. *Orthop. Traumatol. Surg. Res.* **104**, S19–S24 (2018).
22. Westblade, L. F., Simon, M. S. & Satlin, M. J. Bacterial Coinfections in Coronavirus Disease 2019. *Trends Microbiol.* (2021) doi:10.1016/j.tim.2021.03.018.
23. He, S. *et al.* Clinical characteristics of COVID-19 patients with clinically diagnosed bacterial co-infection: A multi-center study. *PLOS ONE* **16**, e0249668 (2021).
24. Lammie, P., Solomon, A., Secor, E. & Peeling, R. *DIAGNOSTIC NEEDS FOR NTD PROGRAMS. The Causes and Impacts of Neglected Tropical and Zoonotic Diseases: Opportunities for Integrated Intervention Strategies* (National Academies Press (US), 2011).
25. Berg, R. D. Bacterial Translocation from the Gastrointestinal Tract. in *Mechanisms in the Pathogenesis of Enteric Diseases 2* (eds. Paul, P. S. & Francis, D. H.) 11–30 (Springer US, 1999). doi:10.1007/978-1-4615-4143-1_2.

26. Vaishnavi, C. Translocation of gut flora and its role in sepsis. *Indian J. Med. Microbiol.* **31**, 334–342 (2013).
27. Fine, R. L., Manfredo Vieira, S., Gilmore, M. S. & Kriegel, M. A. Mechanisms and consequences of gut commensal translocation in chronic diseases. *Gut Microbes* **11**, 217–230 (2020).
28. Bischoff, S. C. *et al.* Intestinal permeability – a new target for disease prevention and therapy. *BMC Gastroenterol.* **14**, 189 (2014).
29. Rendeiro, A. F. *et al.* Longitudinal immune profiling of mild and severe COVID-19 reveals innate and adaptive immune dysfunction and provides an early prediction tool for clinical progression. 2020.09.08.20189092 <https://www.medrxiv.org/content/10.1101/2020.09.08.20189092v1> (2020) doi:10.1101/2020.09.08.20189092.
30. Rendeiro, A. F. *et al.* Metabolic and immune markers for precise monitoring of COVID-19 severity and treatment. 2021.09.05.21263141 <https://www.medrxiv.org/content/10.1101/2021.09.05.21263141v1> (2021) doi:10.1101/2021.09.05.21263141.
31. BBMap. *SourceForge* <https://sourceforge.net/projects/bbmap/>.
32. Krueger, F. & Andrews, S. R. Bismark: a flexible aligner and methylation caller for Bisulfite-Seq applications. *Bioinformatics* **27**, 1571–1572 (2011).
33. Song, Q. *et al.* A Reference Methylome Database and Analysis Pipeline to Facilitate Integrative and Comparative Epigenomics. *PLoS ONE* **8**, e81148 (2013).
34. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
35. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. metaSPAdes: a new versatile metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).
36. Menzel, P., Ng, K. L. & Krogh, A. Fast and sensitive taxonomic classification for metagenomics with Kaiju. *Nat. Commun.* **7**, 11257 (2016).
37. Shoaie, S. *et al.* Global and temporal state of the human gut microbiome in health and disease. (2021) doi:10.21203/rs.3.rs-339282/v1.