

An improved workflow for accurate and robust healthcare environmental surveillance using metagenomics

Jiaxian Shen

Northwestern University <https://orcid.org/0000-0003-4929-8955>

Alexander McFarland

Northwestern University <https://orcid.org/0000-0002-1803-3623>

Ryan Blaustein

National Institutes of Health

Laura Rose

Centers for Disease Control and Prevention

K. Perry-Dow

Centers for Disease Control and Prevention

Mary Hayden

Rush Medical College <https://orcid.org/0000-0002-4603-8501>

Vincent Young

University of Michigan Medical School <https://orcid.org/0000-0003-3687-2364>

Erica Hartmann (✉ erica.hartmann@northwestern.edu)

Northwestern University <https://orcid.org/0000-0002-0966-2014>

Article

Keywords:

Posted Date: February 15th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1303703/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

29 **Introduction**

30 Effective microbial surveillance in the built environment is increasingly important in infection
31 prevention, given the persistence of pathogens in environmental reservoirs and their potential
32 transmission to patients¹⁻⁷ (e.g., carbapenem-resistant *Klebsiella pneumoniae* in sink drains^{8,9}).
33 Furthermore, with the emergent studies showing synergistic relationships among pathogens on
34 hospital surfaces¹⁰ and the possibility for pathogenic bacteria to acquire antibiotic resistance genes
35 from non-pathogenic neighbors¹¹, it is necessary to expand from targeted surveillance to
36 untargeted methods. Untargeted methods are advantageous in identifying novel or rapidly
37 emerging pathogens¹². Metagenomics-based techniques are the most promising option to achieve
38 these goals but are currently challenged by several limitations: 1) they are not powerful enough to
39 extract valid signals out of the background noise for low biomass samples; 2) they do not
40 distinguish between viable and non-viable organisms; and 3) they do not reveal the microbial load
41 quantitatively^{13,14}.

42 For challenge 1), adoption of appropriate negative controls has been emphasized^{15,16}, along with
43 various bioinformatic filtering tools to remove putative contaminants^{17,18}. While current efforts
44 have largely focused on contamination prevention, increasing the biomass itself remains
45 understudied¹². Having adequate biomass is essential, as previous work has indicated that a small
46 amount of starting material (i.e., DNA) has adverse impacts on the outcome regardless of sample
47 processing methods¹⁹. In practice, methods have been adopted as temporary fixes, such as pooling
48 samples from different sites or dates²⁰, and using wipes instead of swabs as sample collectors¹.
49 However, these workarounds are not always available²¹. Moreover, systematic evaluation and
50 benchmarking of optimization strategies for metagenomic sample preparation remain largely
51 unexplored.

52 For challenge 2), propidium monoazide (PMA) is the most widely used viability indicator
53 compatible with molecular techniques. Though intensively optimized²², its efficacy and
54 applicability in combination with metagenomics are controversial. A semi-quantitative systematic
55 evaluation concluded that PMA treatment coupled with 16S rRNA gene amplicon sequencing
56 (PMA-Seq) is reliable when the microbial community is not very complex, while uncertainties
57 increase dramatically with complexity²³. The uncertainties come from both heterogeneity of
58 microorganisms (e.g., cell envelope structure differences, spore formation), and complexity of the
59 background matrix (e.g., turbidity, salt concentration, dead cell density)^{13,24,25}. While the microbial
60 communities to be surveilled have their inherent advantage of being low complexity, little is
61 known about the effectiveness of incorporating PMA with multi-species internal standards. To be
62 appropriately rigorous, comparisons are needed relative to standard surveillance that does not
63 consider viability (i.e., no PMA), as well as traditional methods (i.e., cultivation).

64 For challenge 3), pitfalls of using relative abundances in microbial profiling have been widely
65 indicated. Such pitfalls include but are not limited to lack of unique connections between biological
66 interpretations and experimental observations and unreliable comparisons across samples^{14,26,27}.
67 Strikingly, mis-selection of analytical tools for relative abundance data could lead to as high as
68 100% false discovery rates²⁶. Besides flow cytometry^{28,29}, combing sequencing with quantitative
69 PCR (qPCR) and including internal standards³⁰ are two major means of making quantitative

70 estimations out of next-generation sequences. In practice, previously reported applications for
71 qPCR include air and dust samples in classrooms³¹; for internal standards, applications include
72 Amazon River plume³², soil³³, and stool samples³⁴. However, comparisons are not yet available
73 between metagenomics coupled with qPCR and with internal standards using low biomass
74 environmental samples in the immediate vicinity of humans (e.g., healthcare settings).

75 An additional practical challenge in developing a robust pipeline with metagenomics is how deep
76 one should sequence. While useful in whole genome sequencing, recommendations of coverages
77 expressed by folds of genome sizes (e.g., 15X to 60X) are not readily transferable to metagenomic
78 sequencing (MetaSeq), as reads do not equally distribute across members with substantially
79 different abundances. Nonpareil, a redundancy-based tool, estimates and projects abundance-
80 weighted average coverage for metagenomics (expressed in percentage)³⁵⁻³⁷. This helps reduce
81 erroneous interpretations out of metagenomic results. Yet expected coverage is still largely
82 unpredictable before sequencing is run. Researchers usually rely on previous experience of similar
83 samples and the available budget to determine the sequencing effort (read size, unit: bp), which
84 may lead to either a coverage too low, thus limiting the extractable information¹¹, or a waste of
85 resources¹.

86 To address these challenges, we present a workflow for metagenomics-based environmental
87 surveillance that is appropriate for low-biomass samples, distinguishes viability, is quantitative,
88 and estimates sequence resources (**Fig. 1**). Liquid-liquid extraction, PMA treatment equipped with
89 internal standards and absolute abundance profiling, qPCR, and a machine learning-based model
90 are the recommended components for the comprehensive workflow, with whole-cell filtration and
91 cultivation as optional accessories.

92 **Results**

93 **Liquid-liquid extraction improves the power of handling low-biomass samples**

94 To improve DNA yield of low biomass samples, we first compared 3 categories of extraction
95 methods. Bead-beating and heat lysis followed by liquid-liquid extraction was the optimal method,
96 as opposed to widely used column- and magnetic bead-based methods (“Methods”,
97 **Supplementary Fig. 2**)^{38,39}. Notably, no detectable DNA was recovered using Qiagen DNeasy
98 PowerSoil Kit. Supplementary to the recommendation that DNA input ≥ 1 ng for Nextera Flex
99 Library Prep kit, we correlated it to the practical outcome and found that DNA > 11.2 ng
100 corresponded to raw reads $> 1e+05$ (**Supplementary Fig. 2**).

101 In addition to being low biomass, environmental samples of interest are usually in the immediate
102 vicinity of humans, and thus often contain eukaryotic cells. These cells may compete with bacteria
103 for sequencing depth, lowering detectable resolution on bacteria, especially for low-abundance
104 members. Collection methods such as swabs and wipes can further recover abiotic debris along
105 with biological materials, to which chemicals potentially interfering with downstream experiments
106 may adsorb⁴⁰. To address these issues, we evaluated the implementation of whole-cell filtration in
107 the workflow. Four filtration steps (100, 80, 41, 5 μm) were conducted in descending order of pore

108 sizes. For our samples, filtration did not exert a significant effect on detected proportions of
109 bacteria (**Fig. 2a**) or eukaryotic reads (**Supplementary Fig. 3**), according to paired t tests ($p \geq$
110 0.05). Considering that the non-bacterial proportion of our samples was relatively small (~1%),
111 filtration appears ineffective (or unnecessary) in increasing the bacterial proportion by excluding
112 eukaryotic cells for samples with similar characteristics. Instead, most of the eukaryotic reads were
113 human-associated and thus able to be removed *in silico* (**Supplementary Fig. 3**). Moreover, we
114 did not observe an increase in the number of rare taxa post filtration. Nevertheless, filtration did
115 not negatively affect the number of recoverable taxa (**Supplementary Fig. 4**)^{41,42}.

116 As expected, filtration introduced biomass loss of ~13-44%, according to 16S rRNA gene copy
117 number (**Fig. 2b**). The biomass loss may be compensated by a two-fold concentration, material
118 permitting. Alternatively, the total biomass loss can be reduced in practical applications where
119 one-step filtration is streamlined. Filtration did not impact the overall bacteria composition (**Fig.**
120 **2c**), nor did it change the relative abundances of top abundant taxa (average abundance > 1%)
121 (**Supplementary Fig. 5**). This evidence supports the validity of using filtration to concentrate
122 bacterial samples in sequencing-based experiments for profiling relative abundances. However,
123 the absolute abundances would be affected disproportionately, as the extent of biomass loss varied
124 across samples (**Fig. 2b**). Therefore, when an absolute metric is of interest, the recovery rate needs
125 to be rigorously measured. Bacteria retention profiles on 5 μm filters were similar to those of the
126 liquid samples. However, bacterial members were not proportionally retained by filters of a larger
127 pore size (100, 80, 41 μm) (**Fig. 2d**). Hence treating microbial samples with large pore-size filters
128 may introduce biases, even when relative abundances are used.

129 Taken together, for samples whose non-bacterial proportion is small (e.g., ~1%), it is unnecessary
130 to incorporate filtration to increase the bacterial fraction. Filtration is valid in concentrating
131 samples. However, for low biomass samples which are low in both cell density and quantity,
132 biomass loss outweighs the slight increase of bacterial signal. Instead, switching to a high-yield
133 DNA extraction method, such as liquid-liquid extraction, can achieve higher folds of signal
134 improvement (DNA concentration from undetectable to $18.62 \pm 1.16 \text{ ng}/\mu\text{L}$).

135 **PMA and cultivation improve the ability to determine viability**

136 We examined the efficiency of PMA treatment coupled with metagenomic sequencing (PMA-
137 MetaSeq) on hospital-associated surface samples with the ZymoBIOMICS Microbial Community
138 spike-in as the internal standard. The Zymo community consists of 8 bacterial species and 2 yeasts,
139 which presumably will function more comprehensively and accurately regarding bias correction
140 and quality control than a single-species standard^{13,23,24,43}. Sequencing outcomes were compared
141 with cultivation results for benchmarking, as cultivation is the gold standard for determining
142 microbial viability.

143 Absolute abundance of samples decreased after PMA treatment, indicating the depletion of non-
144 viable signals (**Fig. 3c**). This was further supported by the observation that α diversity was lower
145 in PMA-treated samples (**Fig. 3a**), and that inter distances between paired samples were larger
146 than intra distances within each sample group (Jaccard Distance; **Fig. 3b**). We note that absolute
147 abundance should be used when analyzing sequence data involving viability assessment, as

148 relative abundance profile is likely distorted (**Fig. 3c-d**)³³. Although relative abundance is
149 informative in demonstrating presence/absence, it neglects the amount of overall biomass and thus
150 may inflate the apparent abundance of even low-abundance organisms. While absolute abundance
151 is more reflective of reality, field trials are necessary to determine whether absolute or relative
152 abundance, or either, can be linked to infection or other clinical outcomes.

153 We calculated the efficacy of 8 spike-in bacteria²³. The efficacy should be 1 under ideal conditions,
154 given that the percentage of viable microbes in the Zymo community is negligible
155 (**Supplementary Fig. 6**). The efficacy equaled 1 for all taxa, suggesting that PMA treatment is
156 effective in low biomass samples regardless of taxonomy. This conclusion is partly consistent with
157 our previous evaluation of PMA-Seq where *E. coli* was spiked in²³.

158 Focusing on specific taxa (**Fig. 4b**), we observed occasions of a complete depletion for high-
159 absolute-abundance taxa and retention for low-absolute-abundance taxa, suggesting an effective
160 viability distinction. Relative abundance for some taxa increased after PMA treatment
161 (*g_Pseudomonas*, *s_Pseudomonas psychrophila*, *c_Gammaproteobacteria*, *s_Pseudomonas fragi*,
162 *s_Pseudomonas koreensis*, *k_Bacteria*), while all taxa showed a decrease in absolute abundance.
163 This indicates that PMA treatment may increase the ability to detect taxa with majoritarily viable
164 populations. We did not detect new taxa that were previously undetectable in PMA-treated samples.
165 However, if nonviable microbes are not of interest, treating samples with PMA will improve the
166 detection power for the overall community with comparable sequencing resources (**Fig. 6c**), as it
167 reduced the overall α diversity (**Fig. 3a**).

168 Cross referencing between cultivation and PMA-MetaSeq was greatly impeded by their inherent
169 limitations (e.g., detection limit for MetaSeq, biases with bioinformatics; viable but non-culturable
170 cells for cultivation). Even for PMA-untreated samples, cultivation and MetaSeq only agreed with
171 each other on a small number of taxa (**Fig. 4a**). Among the 3 viable taxa confirmed by cultivation,
172 viability of *s_Pseudomonas fragi* and *s_Pseudomonas stutzeri* was reflected by PMA-MetaSeq,
173 while *s_Pseudomonas fluorescens* became undetectable after PMA treatment. This might imply
174 over-depletion, but could also be because its abundance went below the detection limit of MetaSeq.
175 As indicated by Barbau-Piednoir et al., less-abundant taxa were more likely to be eliminated (to
176 undetectable) by PMA treatment⁴⁴. This is consistent with our observation, as the abundance of
177 *s_Pseudomonas fluorescens* was the smallest among the cultivation-confirmed taxa. Thus for low-
178 abundance taxa, cultivation could serve as a supplement to sequence-based viability assessment
179 techniques, as a small unintentional removal of viable cells may lead to a large presence/absence
180 difference. Moreover, incorporating cultivation can expand the detection spectrum in general, and
181 particularly for low-abundance taxa, due to MetaSeq's restrictions such as detection limit and
182 failure to distinguish closely related taxa.

183 Collectively, we emphasize the importance of using absolute abundance and demonstrate a
184 successful application of multi-species internal standards in PMA-MetaSeq. PMA is effective in
185 low biomass samples and can improve the detection power by eliminating irrelevant signals.
186 Cultivation remains a valuable supplement to sequence-based techniques for capturing a
187 comprehensive picture.

188 **Poor taxonomic classification is a major hurdle for internal standards in quantitative**
189 **metagenomics**

190 Quantifying metagenomics-based abundances using internal standards has substantial benefits.
191 Theoretically, addition of internal standards could compensate for errors resulting from non-
192 quantitative steps¹³. *E. coli* is one of the most used spike-in strains, in part because it is well-studied
193 and easy to recover in sequencing²³. However, ideally, we want the internal standard to contain a
194 set of diverse taxa, so that it well represents the diversity in microbial communities. We
195 investigated the performance of the Zymo community as the internal standard for hospital-
196 associated environmental samples, along with qPCR for the 16S rRNA gene.

197 Unfortunately, the efficiency of implementing the Zymo standard in quantitative metagenomics
198 was drastically impeded by the limited resolution of taxonomic classification. We tried two
199 approaches: Metaxa2⁴⁵⁻⁴⁷ coupled with the SILVA 132 SSU database^{48,49} and MetaPhlAn3^{50,51},
200 which uses a collection of marker genes. The taxonomic resolution varied substantially across
201 different taxa. For samples containing only the Zymo standard, 85% of the small subunit rRNA
202 reads were attributable by Metaxa2, while only 48.14% of the metagenomes were recognized by
203 MetaPhlAn3. Within the attributable portion, MetaPhlAn3 performed better regarding specificity;
204 all reads were classified at the species level, while Metaxa2 retained a decent amount of
205 information at higher levels, with the ratio of genus/family-level and species-level classifications
206 ranging from 0.18 to 11.28.

207 Foreseeably, this issue will be alleviated as reference databases and taxonomic assignment tools
208 continue to advance. Currently, advantages of internal standards are mainly reflected when
209 species-level identification is the major focus. For instance, clinical samples usually target
210 pathogenic species whose core pangenomes are relatively well represented in databases. In this
211 case, the biases from uneven representation of species can also be corrected based on the
212 performance of closely related internal-standard species. However, if information at genus or
213 higher levels is of consideration, internal-standard techniques become non-applicable, as we are
214 not able to distinguish internal-standard taxa from other species within the same genus (or at higher
215 levels), which is the basis of making calculations and corrections. Coupling with qPCR, instead,
216 is more appropriate (**Fig. 3c**). Environmental communities are typical examples where coupling
217 with qPCR stands out because environmental microorganisms are not usually well represented at
218 the species level. Of 87 samples in our study, strikingly, MetaPhlAn3 only recognized an average
219 of 19.68% of the metagenome at the species level. The classification rate slightly increased to
220 38.40% using Metaxa2, which substantially improved to 87.24% when genus level was included.

221 **Accessible sample features can predict required sequencing effort**

222 To enable more informed decision-making before MetaSeq, we conducted a quasi-meta-analysis,
223 using the limited number of existing hospital-related environmental metagenomics studies^{1,5,20,52-}
224 ⁵⁴. We recruited 956 shotgun samples (874 from 6 previous studies and 82 from this study)
225 (**Supplementary Table 1**). Using these data, we linked accessible features (e.g., location, building,
226 sampling method) to the required sequencing effort given a targeted coverage, leveraging machine
227 learning-based models and Nonpareil (**Fig. 6a**).

228 Relationships between Nonpareil diversity (Nd, unit: log-bp) and metadata features were first
229 explored (**Fig. 6a: stage one**). Nd is an index measuring the complexity of a microbial community
230 regarding "sequence space", which correlates with classic bin-based diversity indices (e.g.,
231 Shannon index) for bacteria^{35,37}. Though not passing the normality test (Shapiro-Wilk Test, $p =$
232 $3.338e-16$)⁵⁵, normal distribution was still the best-fit distribution of our dataset, followed by
233 logistic distribution, upon investigation by Cullen and Frey Graph and R package "fitdistrplus"
234 (**Fig. 5a, Supplementary Fig. 7**). Presumably, the deviation from normality will decrease as
235 sample sizes increase. For 90% of samples, Nd was within 2 orders of magnitude (15.4-20.0,
236 natural log scale), suggesting a common range for hospital-associated environmental samples,
237 which is valuable for reference when designing future studies. Notably, this Nd level was among
238 the lowest across 6 different environments including animal hosts, fresh water, and soil (**Fig. 5b**)³⁷.

239 We further examined the influences of sample type (sink versus surface), sampling method and
240 sample pooling on Nd. No significant difference was observed between sink and surface samples
241 (**Fig. 5c**). Within sink samples, Nd was significantly different across sampling methods (Anova, $p =$
242 $7.18e-06$). Specifically, samples collected by swabs seemed to have a smaller diversity than
243 those by the other methods (Tukey's post-hoc test; samples without a clear collection method
244 stated in the original paper were assigned as "Sink"). Note that even though sink samples are
245 generally from the same location, the confounding effects introduced by sub-locations (e.g., sink
246 basin, pipe edge, p trap) cannot be ruled out. Similarly, within surface samples, though Nd of wipes
247 was significantly larger than that of swabs (unpaired t-test), confounding effects remain (e.g.,
248 researchers tend to use wipes for large-area and high-biomass locations, like floors, which often
249 contain more diverse communities) (**Fig. 5d**). Though weak, we noticed a trend of diversity
250 increase after sample pooling (**Fig. 5e**), raising the alarm that more caution should be taken when
251 increasing biomass by pooling samples. The practice of sample pooling assumes that the pooled
252 samples share some core features, whose biomass will be increased past the detection limit. This
253 may be true of certain sample types, e.g., host-associated microbiomes, but is unlikely to be true
254 of built-environment samples that lack a conserved core⁵⁶. Further investigations are needed
255 should more data become available, as the sample size was quite limited for some groups (e.g.,
256 $n(\text{pooled monitor}) = 2$, $n(\text{not-pooled monitor}) = 4$). In the interim, we recommend seeking other
257 methods, such as a high-yield DNA extraction, before resorting to sample pooling, as the resulting
258 sample characteristics may be different from individual samples.

259 To further harness the reference potential of Nd, we built models to predict Nd from metadata
260 features based on machine learning algorithms. Eight predictor variables (location, building, study,
261 country, touch frequency, sample type, sampling method, sample pooling) were included based on
262 data availability, MIxS-BE standards, and previous experience (**Supplementary Table 1**)^{53,57,58}.
263 Nd, the response variable, was first converted from a numerical variable to a nominal variable.
264 Three conversion schemes were tried, with the intervals being 2.5, 1.0, and 0.5 (number of
265 categories being 2, 5, 11, respectively). Random sampling was adopted to split the entire dataset
266 into training and testing datasets at the ratio of 4:1. Implementing repeated cross-validation (5
267 folds, 5 times) on the training dataset, 9 algorithms were examined to optimize the classification
268 performance, including random forest, stochastic gradient boosting, and support vector machines.
269 Algorithms were evaluated according to 4 metrics (area under curve, Kappa, and balanced

270 accuracy on both training and testing datasets)⁵⁹. Overall, no difference was observed among the
271 tested algorithms. Random forest was selected due to its slightly better performance from a holistic
272 perspective and capability of ranking the predictor variables.

273 The model accuracy positively correlated with the interval size. At 2.5, the accuracy on the training
274 dataset was as high as 87.69%, and slightly lower on the testing dataset (82.60%). The accuracy
275 dropped as the classification demand rose. The mean balanced accuracy on the testing dataset was
276 64.08% and 61.38% when intervals were 1.0 and 0.5, respectively. Considering that Nd was
277 converted from a continuous variable, we examined the misclassifications and found that most of
278 them fell into nearby categories. We thus calculated the mean balanced accuracy ± 1 category and
279 observed a substantial improvement. Specifically, mean balanced accuracy of 87.06% and 77.17%
280 can be achieved for 5- and 11-category classifications, respectively. Considering the sparsity of
281 the currently available dataset and the challenge of multiclass classification, this model
282 demonstrated a reasonable degree of accuracy, which should improve as sample sizes and available
283 features grow.

284 The variable importance ranking generated by random forest separated the predictor variables into
285 3 groups (**Fig. 6b**). Location, building, and study were the top 3 variables with the highest
286 importance, followed by country, touch frequency, and sample type, while sample pooling and
287 sampling method hardly impacted the classification. Group-wise, this ranking was generally
288 consistent with the explanatory power described by linear regression (**Supplementary Fig. 8**)⁵⁸.
289 That “study” ranked as one of the most important variables indicated the existence of biases
290 towards individual studies in the current dataset (e.g., “batch effects” related to respective
291 sampling, processing, sequencing, and analysis), which was also observed by a previous meta-
292 analysis of indoor microbiota⁵⁸. Interestingly, despite its high importance, the model performance
293 had almost no drop after excluding “study” (> 95% for all conversion schemes), justifying making
294 predictions without involving artificial metadata features like “study”. It is worth noting that
295 importance of the other variables (e.g., building, country) was raised after this exclusion
296 (**Supplementary Fig. 9**). To find the features necessary for making a comparably accurate
297 prediction, we further examined the performance of models after gradually reducing the number
298 of predictor variables, and found that using only 2, “location” and “building”, the new model
299 achieved 95% accuracy regardless of interval sizes tested.

300 With Nd and metadata features connected, the required sequencing effort at a targeted coverage
301 was then inferred (**Fig. 6a: stage two**). Upon fitting the data, we revealed a linear relationship
302 between the natural log of estimated sequencing effort at 95% coverage ($\ln(\text{LRstar})$) and Nd, with
303 the equation being $\ln(\text{LRstar})=1.14*\text{Nd}+1.21$ (Adjusted R-squared = 0.6012, $p < 2.2e-16$) (**Fig.**
304 **6c**). This is theoretically backed up by previous findings that sequencing effort depends on the
305 diversity level and the genome size, and that the latter can be ignored for most microbial
306 communities, particularly bacterial communities, since the differences in genome size are usually
307 no more than one order of magnitude³⁵. Instructions to make calculations between sequencing
308 effort and other coverage levels are provided at
309 https://github.com/jxshen311/workflow_metagenomic_environmental_surveillance/tree/main/nopareil/example_SeqEffort%26Coverage.
310

311 **Discussion**

312 Although sequence-based environmental surveillance of microbial communities for better
313 management of public health has been appealed for and utilized, best practices of the workflow
314 have not been systematically studied to ensure proper interpretations of sequencing results to aid
315 in infection risk assessment^{13,14}. This study introduces a well-structured and informed
316 metagenomics-based workflow towards the goal of being appropriate for low-biomass, viability,
317 quantification, and resource estimation. We recommend adopting liquid-liquid extraction to
318 improve DNA yield and only incorporating whole-cell filtration when non-bacterial proportion is
319 large. Despite its imperfection, we suggest including PMA treatment, and involving cultivation
320 when demanding comprehensive profiling. We further recommend integrating internal standards
321 for quantification, and additionally qPCR when we expect poor taxonomic classification. We also
322 introduce a machine learning-based model to predict required sequencing effort from accessible
323 sample features. The model helps make full use of sequencing resources and achieve desired
324 outcomes.

325 While using realistic samples in testing simulates conditions the workflow may face in practical
326 applications, it comes with side effects. Our aggregation sample had a small fraction of non-
327 bacterial organisms (~1%). Thus, the conclusion that whole-cell filtration does not increase the
328 bacterial proportion and signal of rare taxa to a statistically significant degree is probably only
329 applicable to samples with similar characteristics, representing 84.90% among the 874 samples
330 from hospital-related environmental studies used in the quasi-meta-analysis^{1,5,20,52-54}. However, a
331 few samples did contain a decent proportion of eukaryotes. Specifically, 132 samples harbored
332 more than 1% eukaryotic reads, and strikingly more than half reads were attributed to eukaryotes
333 for 20 samples. Moreover, samples collected from high-touch surfaces were more likely to have
334 higher proportions of eukaryotes than low-touch surfaces and sinks. Of the 104 sink samples, the
335 maximum percentage of eukaryotes was 0.1%. Therefore, filtration is probably unnecessary for
336 most environmental samples (especially sink samples) and may be beneficial for part of high-touch
337 surface samples (**Supplementary Fig. 10**).

338 Despite being semi-quantitative and entailing considerable uncertainty, involving PMA takes us a
339 step closer to understanding viability, particularly for low biomass samples whose complexity is
340 also relatively low²³. Notably, the overall uncertainty comes not only from PMA treatment but
341 also from the metagenomics pipeline itself, such as biases from DNA extraction kits and taxonomic
342 assignment tools^{16,60}. Integration with other omics techniques (e.g., metatranscriptomics,
343 metaproteomics) was proposed to make up for PMA's shortcomings²³. Pursuing viability profiles
344 using orthogonal methods would plausibly enable a more comprehensive understanding, but the
345 cost-benefit ratio may be considerably high for multi-omic techniques. Integrating with cultivation,
346 instead, provides an affordable alternative. Notwithstanding, it remains to be investigated how to
347 properly interpret results generated by a combination of methods, as inconsistencies between
348 disparate methods are common.

349 Some overlap exists between methods for viability determination and those for depleting
350 eukaryotic DNA. For example, osmotic lysis followed by PMA treatment is recommended to
351 remove human DNA in saliva samples⁶¹. However, recommended methods depend on the sample

352 type. PMA is not recommended for sputum samples, where nuclease-based methods (e.g., digest
353 with benzonase) showed an equal or better performance⁶². Factors impacting method performance
354 include percentage and composition (e.g., extracellular DNA, DNA in largely lysed or partially
355 compromised cells) of targets to be removed (i.e., eukaryotes and dead bacteria), as well as
356 characteristics of background matrix (e.g., viscosity). For instance, saliva and sputum consistently
357 contain $\geq 90\%$ human DNA^{61,62}, while this percentage is very diverse for hospital-associated
358 environmental samples (**Supplementary Fig. 10**). Filtration failed to exclude human DNA in
359 saliva likely because extracellular DNA was the dominant component rather than cells⁶¹. For
360 sputum samples where cells are lysed and DNA is no longer protected, nucleases might be quite
361 effective in depleting extracellular DNA, whereas PMA efficacy could be hindered by the viscosity
362 of the matrix⁶². In contrast, in environments where cells gradually decay due to harsh conditions
363 (e.g., desiccation), more DNA attributable to dead cells would still have a partially compromised
364 membrane; PMA, as a small molecule, may be more effective in penetrating the damaged cell
365 membrane and depleting the DNA. For eukaryotic depletion, it may be beneficial to further unravel
366 the underlying mechanisms influencing the efficacy of different methods in different sample types
367 and characteristics. For viability assessment, instead of focusing on this viable/dead dichotomy,
368 perhaps more critically, we should keep in mind that “viability” is rather an intermediate or
369 methodological term, linking surveillance results to questions of interest (e.g., which bacteria are
370 infectious)¹³.

371 We applied multi-taxa internal standards and calculated PMA efficacy of spike-in taxa based on a
372 reasonable assumption that the percentage of viable microbes in the Zymo community is negligible,
373 resulting in a theoretical value of 1. While this internal standard can strongly reflect incomplete
374 suppression of non-viable signals, potential toxicity of PMA might be underrepresented. Although
375 no toxicity was observed at the PMA dose of our protocol in validation (**Supplementary Fig. 11**),
376 a customized internal-standard mixture featuring 0.5 as the designed PMA efficacy would be ideal
377 for future studies²³. As opposed to purchasing commercial products, we recommend utilizing the
378 Zymo community as a reference for the taxonomic composition and constructing the mixture with
379 live cultures in real-time, because viability (or membrane integrity when PMA is used) is difficult
380 to maintain in manufacturing, shipping, and storage.

381 Continuous advancement of internal standards for quality control, as well as quantification and
382 other features, is still one of the major hotspots in method optimization. A suitable internal standard
383 should well balance representation and recognizability. Good representation means that the
384 workflow impacts the spike-in and targeted microbes comparably (because of their similarity).
385 Good recognizability means that the spike-in can be easily distinguished from the targets. In this
386 study, the Zymo community was selected due in large part to its representation, as it spans broadly
387 the phylogenetic tree. Previous studies have selected internal standards based on a similar principle.
388 For instance, the Zymo community and a 10-species mock community were chosen for
389 gastrointestinal and stool samples, respectively^{63,64}. Peroxide-killed *Campylobacter sputorum* was
390 used to quantify viable thermotolerant *Campylobacter*⁶⁵. These internal standards are prone to be
391 confounded with targets, thus posing challenges for bioinformatics to accurately identify and
392 quantify taxa. To obtain good recognizability, exotic materials are sought. In the aforementioned
393 example, the researchers chose 10 species that were generally absent from the stool of healthy

394 individuals. The same criterion was followed by another gut microbiome study in which microbes
395 from hypersaline environments, soil, and plants were utilized³⁴, as well as a study on Amazon
396 River plume to which genomic DNA from *Thermus thermophilus* HB8 was applied³². Finding a
397 completely exotic species is more challenging for environmental surveillance whose subjects are
398 influenced by both human and environmental activities. As a potential solution, artificial DNA
399 have been developed to ensure differentiation from the targets. Previous reports included sets of
400 synthetic DNA, 16S rRNA genes, and chimeric DNA fragments, implemented in different venues
401 of metagenomic and amplicon sequencing^{33,66,67}. However, whether these exogenous (or even
402 artificial) standards' behavior resembles that of the targets remains questionable. By and large,
403 more systematic evaluation and optimization are needed to foster the development of internal-
404 standard techniques that better balance representation and recognizability, or at least make their
405 pros and cons quantitatively accessible, both in general and for specific contexts.

406 Though the classification models performed well from a practical perspective, their accuracy with
407 small intervals still merits improvement. Building a hierarchical classification model might be
408 beneficial, as we observed a drastic increase in the accuracy when the interval size was enlarged.
409 It is also likely that the available dataset is not good enough to train a model with very high
410 accuracy. For example, there is clear evidence that the data were biased by the disparate sample
411 sizes between studies. Moreover, we only managed to collect 7 common metadata features
412 (excluding “study”) without involving a substantial number of missing values, which raises the
413 question of whether what we achieved has already reached the theoretical plateau of explanatory
414 power of these features. If this is the case, standardized reporting of more high-quality metadata
415 should be further promoted. Additionally, since normal distribution was the best-fit distribution of
416 the current dataset, with seemingly missing pieces in the middle (**Fig. 5a**), fitting data into known
417 distributions may be more explanatory as large sizes of data become accessible.

418 This metagenomics-based environmental surveillance workflow is particularly useful in infection
419 prevention and disinfection assessment. Although we focus on microbial surveillance of built
420 environments, especially hospital-associated surfaces, the workflow developed in this study can
421 be adapted to other contexts with similar characteristics. For example, the multifaceted lessons
422 learned from this study will benefit the continuing development of microbiome-based clinical
423 testings from body sites (e.g., skin), such as methods to increase low-biomass signals and
424 determine viability¹². Moreover, the experience gained in overcoming challenges unique to
425 environmental microbiomes (e.g., quantitative metagenomics with poor taxonomic classifications)
426 is also useful to studies on other environments, such as wastewater and air.

427 **Methods**

428 **Sample collection, aggregation, and cultivation**

429 We collected 120 surface swabs from the 28-bed medical intensive care unit (MICU) at Rush
430 University Medical Center (RUMC) in October 2018. RUMC is a 720-bed tertiary care teaching
431 hospital in Chicago, IL. Samples were collected from door sills, computer keyboards, light
432 switches, nurse call buttons, and bed rails in 13 single-bed patient rooms, as well as door sills in 4

433 medication rooms, 2 public restrooms, 1 staff-only restroom, and the communicating space of
434 MICU (**Supplementary Table 2**). Weighted mean area of sampled surfaces was 216 cm². Patient
435 rooms were selected to keep a relatively balanced number for both contact isolation and non-
436 contact isolation rooms. Healthcare providers and visitors entering contact isolation rooms are
437 required to wear gowns and gloves, which may reduce transmissions via contaminated healthcare
438 providers. Room temperature and relative humidity were documented during the collection, which
439 varied slightly across rooms, with the average being 23.8°C and 45%, respectively. Each sample
440 was collected by 3 COPAN Nylon Flocked Swabs (Copan Diagnostics, Murrieta, CA, USA) and
441 1.5 mL Phosphate Buffered Saline with 0.02% Tween 80 (PBST), and stored at 4°C for up to 24
442 h prior to extraction, aggregation, and cultivation^{53,68}. Swabs were extracted and aggregated to
443 create a representative microbiome sample^{56,68,69}. Aliquots of this aggregation sample were then
444 subjected to different processing methods (i.e., several DNA extraction methods, microbial
445 community standard spike-in, PMA treatment and whole-cell filtration) to find best practices of
446 the workflow (**Supplementary Fig. 1**).

447 To capture a large fraction of the indoor microbiome diversity, we cultured the samples with 4
448 different media: tryptic soy agar (TSA), Reasoner's 2A agar (R2A), 0.1 strength R2A at 25°C, and
449 blood agar (BA) at 37°C, all supplemented with 4 mg/L itraconazole⁶⁸. This resulted in 233
450 cultivable isolates. All colonies that could be individually picked or purified were subject to
451 taxonomic identification by matrix-assisted laser desorption/ionization time-of-flight mass
452 spectrometry (MALDI-TOF MS) using the VITEK® MS Mass spectrometry microbial
453 identification system (BioMerieux, Marcy-l'Étoile, France) and the VITEK MS V3.2 FDA 510(k)
454 cleared database. Among the 233 isolates, 201 were identified. It is important to note that, because
455 multiple media types were used, the number of isolates for each species identified does not
456 represent the relative abundance of this species in the sample, as some species may have grown on
457 multiple media.

458 **Standard addition, PMA treatment and whole-cell filtration**

459 All treatments were done in triplicate, including cultivation.

460 Standard addition

461 Aliquots were snap frozen and stored at -80 °C until further processing to maximize the integrity
462 of samples and avoid degradation resulting from long-term storage at 4°C^{69,70}. Samples were
463 thawed at 4°C prior to treatments. ZymoBIOMICS Microbial Community Standard (Zymo
464 Research, Irvine, CA, USA) was used as both the internal standard and the external standard. As
465 the internal standard, 6.50 µL Zymo community was spiked into 1 mL aggregate sample, following
466 the criterion that DNA of the species with the highest abundance in the Zymo community
467 approximates 1% of the total DNA in the aggregate sample^{32,71}. As the external standard, aliquots
468 of the Zymo community were run in parallel with aggregate samples throughout the workflow to
469 assure its performance.

470 PMA treatment

471 Following standard addition, PMA treatment (Biotium, Fremont, CA, USA) with an optimized
472 protocol was applied to half of the samples within each group^{22,24,25,72,73}. The protocol was first
473 validated by reproducing the work of Nocker et al. (2006) using *Escherichia coli* (ATCC 8739) as
474 model strain. *E. coli* was grown to the exponential phase and half killed by heat inactivation at
475 95°C for 7 min in Eppendorf ThermoMixer shaking at 400 rpm for homogenized heating. After
476 cooling to room temperature, live and heat-killed cultures were mixed following the same ratios
477 as in Nocker et al. (2006), yielding samples of 6 different expected live cell ratios. Viability of
478 both live and heat-killed cultures was confirmed by spread plating onto TSA and incubating at
479 37°C overnight. Half of the constructed samples underwent PMA treatment. The results were
480 evaluated by both DNA concentration ratio quantified by Quant-iT™ PicoGreen™ dsDNA Assay
481 (ThermoFisher, Waltham, MA, USA) and copy number ratio by qPCR with 16S universal primers
482 (341F and 534R) (**Supplementary Fig. 11**). Briefly, a final concentration of 25 µM PMA was
483 used, and several steps were conducted to ensure the consistency across samples and minimize
484 nonspecific reactions between PMA and random sample components, including 1) adding PMA
485 to tube caps and inverting all tubes simultaneously, 2) working under red light, and 3) protecting
486 samples from light as much as possible before the light activation step in the PMA-Lite™ device
487 (Biotium, Fremont, CA, USA). An aliquot of samples for each replicate was preserved at -80°C
488 until DNA extraction, with the rest stored at 4°C for downstream filtration.

489 Whole-cell filtration

490 Whole-cell filtration was conducted using EMD Millipore 25 mm Glass Vacuum Filter kit
491 (MilliporeSigma, Burlington, MA, USA), 125 mL filter flask, and Gemini vacuum pump in a
492 biosafety cabinet following aseptic techniques (**Supplementary Fig. 12**). Notably, autoclaved
493 tweezers were used to avoid possible contaminations from touching sensitive parts of the set-up.
494 Samples were filtered by 100 µm nylon membrane, followed by 80 µm and 41 µm nylon
495 membranes and 5 µm PVDF membrane (MilliporeSigma, Burlington, MA, USA). 1 mL PBS was
496 added to the falcon tube and flask at each step to increase the sample recovery by rinsing the inner
497 wall. The filtered samples were then subjected to 3-fold (relative to the volume before filtration)
498 vacuum concentration with an Eppendorf Vacufuge plus. Filtered liquid samples and filter papers
499 were preserved at -80°C until DNA extraction. To increase the extraction efficiency from filter
500 papers, we compared 1) cutting them with scissors into 9 pieces, 2) grinding them with metal
501 spatula after snap freezing in liquid nitrogen, and 3) directly putting the whole filter paper into the
502 preservation tube. We finally selected the third option as this was the most operationally feasible
503 way without high risk of contamination.

504 **Negative controls**

505 To combat the susceptibility of low-biomass samples to contamination, we included 4 types of
506 negative controls along the workflow, namely, 6 negative field controls, 6 negative media controls,
507 12 negative filter controls, and 7 negative kit controls^{13,15,16,58}. The negative controls were
508 processed in parallel with the surface samples, including metagenomic sequencing and
509 bioinformatic analysis.

510 Negative field controls were collected exactly the same as surface samples, except that the swabs
511 were exposed to the air without contacting targeted surfaces. One negative field control was
512 collected at the beginning and the end of each sampling session. Two negative media controls were
513 included in each sampling session, which were unopened media with swabs from the same lot.
514 Each collected control was split into triplicate and processed along with samples^{53,74}. Negative
515 filter controls were included in triplicate for each pore size by letting sterile PBST flow through
516 the vacuum filtration system attached with blank filter papers. Additionally, 7 negative kit controls
517 were processed across batches of DNA extractions.

518 **DNA extraction, qPCR and metagenomic sequencing**

519 To ensure enough DNA recovery, we performed an initial optimization on a separate set of surface
520 swab samples collected from the same MICU prior to working with the aggregate sample
521 (**Supplementary Fig. 2**). Column-based methods were first tried due to its widespread usage in
522 the field. We examined Qiagen DNeasy PowerSoil Kit with standard protocol and a modified
523 version by 1) changing from vortex lysis to bead-beating lysis, 2) introducing heat incubation after
524 bead-beating, and 3) adding 50 μL water each time for twice in total at the elution step. DNA yields
525 of both were below the limit of detection. Liquid-liquid extractions were performed afterwards for
526 their high-yielding potentials. Phenol-chloroform extraction resulted in the highest yield
527 (186.27 ± 55.51 ng/ μL by NanoDrop), but the purity indicated by 260/280 was not acceptable
528 (1.36 ± 0.03). Lucigen MasterPureTM Complete DNA and RNA Purification Kit also resulted in
529 high yields when coupled with bead-beating and heat lysis and better purity than phenol-
530 chloroform extraction (260/280 1.62 ± 0.02). We attempted to improve the purity using the
531 Agencourt AMPure XP PCR Purification kit. However, we did not see a purity increase (260/280
532 1.62 ± 0.03) and incurred a 64.70% DNA yield drop. Based on the above tests, we noticed that
533 methods involving columns (Qiagen PowerSoil) or magnetic beads (Agencourt AMPure) greatly
534 decreased the DNA yield. Because the primary concern for surface samples is low biomass,
535 increasing DNA yield is considered more critical than bringing 260/280 to the desired range of
536 1.8-2.0. Therefore, the Lucigen MasterPureTM Complete DNA and RNA Purification Kit with the
537 adapted protocol was chosen for all subsequent analyses⁷⁵. Samples were thawed at 4°C prior to
538 DNA extraction and DNA concentrations were quantified by Quant-iTTM PicoGreenTM dsDNA
539 Assay⁷⁶.

540 V3 region of the 16S rRNA gene was amplified in qPCR using universal primers (341F: 5'-CCT
541 ACG GGA GGC AGC AG-3', 543R: 5'-ATT ACC GCG GCT GCT GGC A-3')⁴⁰. The 20 μL
542 reaction mixture consisted of 10 μL PowerUpTM SYBRTM Green Master Mix (Applied
543 Biosystems), 0.6 μL forward primer (10 μM), 0.6 μL reverse primer (10 μM), 5.0 μL DNA
544 templates (pre-diluted if necessary), and 3.8 μL nuclease-free water. The reaction was run in
545 technical triplicate on a QuantStudio 3 Real-Time PCR System (Applied Biosystems) with an
546 initial denaturation step at 95°C for 2 min, followed by 40 amplification cycles (95°C, 15 s;
547 56°C, 15 s; 72°C, 1 min) and a melting curve stage (95°C, 15 s; 60°C, 1 min; 95°C, 15 s). No-
548 template control and 5~8 standards were included in each batch to generate the standard curve
549 (efficiency > 90%; $R^2 > 0.99$). Plasmid DNA constructed by TOPOTM TA CloningTM Kit
550 (Invitrogen, Waltham, MA, USA) was used as standards.

551 Extracted DNA was shipped on dry ice to the UMICH Microbiome Core (Ann Arbor, MI, USA)
552 for library preparation using Nextera™ DNA Flex Library Prep Kit and paired-end 250-bp shotgun
553 metagenomic sequencing on an Illumina MiSeq platform (MiSeq Reagent Kit v2). Libraries were
554 normalized before sequencing, and for samples without enough DNA (e.g., negative controls), all
555 available materials were used.

556 **Data analysis**

557 Sequence data processing

558 KneadData (v0.6.1) was first used to clean the shotgun sequences with default parameters. Reads
559 present in the human reference database (hg37_and_human_contamination) and negative controls
560 were filtered out. Metaxa2 (v2.2)⁴⁵⁻⁴⁷ coupled with SILVA 132 SSU database^{49,77} was chosen to
561 generate taxonomic profiles after comparing it with MetaPhlAn2 (v2.6.0)⁷⁸ and MetaPhlAn3
562 (v3.0.7)^{50,51}. The evaluation was conducted based on their performance on external standards and
563 cross validation with cultivation results for untreated aggregate samples. Default parameters were
564 used for all three tools. MetaPhlAn3 was ruled out mainly because it only generates marker genes
565 at the species level and an average of 80.32% metagenome was deemed unknown for our samples.
566 For external standards, both Metaxa2 and MetaPhlAn2 recognized all 8 bacteria species
567 demonstrated in the theoretical composition, but MetaPhlAn2 failed to classify the 2 eukaryotic
568 species. Moreover, it did not classify *Pseudomonas fluorescens* and barely classified *Pseudomonas*
569 *stutzeri* from aggregate samples, while Metaxa2 recognized both. Though Metaxa2 included a few
570 spurious taxa, all can be eliminated by removing singletons. Since the primary goal of this study
571 was to compare techniques and recommend best practices, sensitivity outweighed specificity.
572 Therefore, Metaxa2 was selected and singletons were removed for downstream analyses. Taxa
573 were labeled to the lowest classifiable level, with species level as the ultimate target²³.
574 Metagenomic sequencing coverage for all samples was estimated by Nonpareil (v3.303) under
575 kmer mode using default settings^{35,37}.

576 Statistical analysis

577 Statistical analyses and data visualization were conducted in R (v4.0.4)⁷⁹ with packages such as
578 Nonpareil, vegan, ape, ggplot2⁸⁰, and dplyr. Principal coordinate analysis (PCoA) based on
579 Jaccard metric was performed to demonstrate beta diversity⁵⁸. Differences between groups were
580 determined by Student's T-test or ANOVA coupled with Tukey's post-hoc test, depending on the
581 number of groups under comparison. $P \leq 0.05$ was defined as statistically significant.
582 Significance codes are: $p > 0.05$ (ns), $0.01 < p \leq 0.05$ (*), $0.001 < p \leq 0.01$ (**), $p \leq 0.001$ (***).
583 Package “fitdistrplus” was implemented to find the best-fit distribution for Nonpareil diversity.
584 Machine learning models were trained using the package “caret”.

585 **Data availability**

586 The raw shotgun metagenomic sequencing data are available in the NCBI SRA repository under
587 Bioproject number PRJNA765404. Source code and supplementary data for reproducing analyses

588 are available under MIT license at
589 https://github.com/jxshen311/workflow_metagenomic_environmental_surveillance. Protocols are
590 available at Protocol Exchange with DOIs: 10.21203/rs.3.pex-1656/v1 (sample collection,
591 extraction, and cultivation), 10.21203/rs.3.pex-1657/v1 (snap freezing), 10.21203/rs.3.pex-
592 1659/v1 (PMA treatment), and 10.21203/rs.3.pex-1658/v1 (DNA extraction).

593 **Acknowledgments**

594 This work was supported by the Centers for Disease Control and Prevention (BAA FY2018-
595 OADS-01 Contract 02915). This research was supported in part through the computational
596 resources and staff contributions provided for the Quest high performance computing facility at
597 Northwestern University which is jointly supported by the Office of the Provost, the Office for
598 Research, and Northwestern University Information Technology. We thank Thelma Dangana and
599 Khaled Aboushaala for their help in collecting samples, and the clinical staff of the medical
600 intensive care unit at Rush University Medical Center for their cooperation. We are also grateful
601 to Pamela B Bell and Rachel Beers for their contributions in performing MALDI-TOF MS.

602 **Author's contributions**

603 JS and EMH designed the study. JS conducted the experiments, performed the analyses, and wrote
604 the manuscript. JS and AM collected the samples. AM and RAB provided support in
605 bioinformatics. LJR, AP, MKH, VBY, and EMH supervised the project. All authors have read,
606 edited, and approved of the final manuscript.

607 **References**

- 608 1 Brooks, B. *et al.* Strain-resolved analysis of hospital rooms and infants reveals overlap
609 between the human and room microbiome. *Nature Communications* **8**, 1814,
610 doi:10.1038/s41467-017-02018-w (2017).
- 611 2 Raveh-Sadka, T. *et al.* Evidence for persistent and shared bacterial strains against a
612 background of largely unique gut colonization in hospitalized premature infants. *Isme j* **10**,
613 2817-2830, doi:10.1038/ismej.2016.83 (2016).
- 614 3 Lax, S. *et al.* Longitudinal analysis of microbial interaction between humans and the indoor
615 environment. *Science* **345**, 1048, doi:10.1126/science.1254529 (2014).
- 616 4 Vickery, K. *et al.* Presence of biofilm containing viable multiresistant organisms despite
617 terminal cleaning on clinical surfaces in an intensive care unit. *Journal of Hospital*
618 *Infection* **80**, 52-55, doi:<https://doi.org/10.1016/j.jhin.2011.07.007> (2012).
- 619 5 Constantinides, B. *et al.* Genomic surveillance of *Escherichia coli* and *Klebsiella* spp. in
620 hospital sink drains and patients. *Microb Genom* **6**, doi:10.1099/mgen.0.000391 (2020).
- 621 6 Martineau, C. *et al.* *Serratia marcescens* Outbreak in a Neonatal Intensive Care Unit: New
622 Insights from Next-Generation Sequencing Applications. *Journal of Clinical Microbiology*
623 **56**, e00235-00218, doi:doi:10.1128/JCM.00235-18 (2018).

624 7 Hu, H. *et al.* Intensive care unit environmental surfaces are contaminated by multidrug-
625 resistant bacteria in **biofilms**: combined results of conventional culture, pyrosequencing,
626 scanning electron microscopy, and confocal laser microscopy. *J Hosp Infect* **91**, 35-44,
627 doi:10.1016/j.jhin.2015.05.016 (2015).

628 8 Burgos-Garay, M. *et al.* Colonization of carbapenem-resistant *Klebsiella pneumoniae* in a
629 sink-drain model biofilm system. *Infect Control Hosp Epidemiol* **42**, 722-730,
630 doi:10.1017/ice.2020.1287 (2021).

631 9 Kotay, S., Chai, W., Guilford, W., Barry, K. & Mathers, A. J. Spread from the Sink to the
632 Patient: In Situ Study Using Green Fluorescent Protein (GFP)-Expressing *Escherichia coli*
633 To Model Bacterial Dispersion from Hand-Washing Sink-Trap Reservoirs. *Applied and
634 environmental microbiology* **83**, e03327-03316, doi:10.1128/AEM.03327-16 (2017).

635 10 D'Souza, A. W. *et al.* Spatiotemporal dynamics of multidrug resistant bacteria on intensive
636 care unit surfaces. *Nat Commun* **10**, 4569, doi:10.1038/s41467-019-12563-1 (2019).

637 11 Ben Maamar, S. *et al.* Mobilizable antibiotic resistance genes are present in dust microbial
638 communities. *PLoS Pathog* **16**, e1008211, doi:10.1371/journal.ppat.1008211 (2020).

639 12 Selway, C. A., Eisenhofer, R. & Weyrich, L. S. Microbiome applications for pathology:
640 challenges of low microbial biomass samples during diagnostic testing. *The Journal of
641 Pathology: Clinical Research* **6**, 97-106, doi:10.1002/cjp2.151 (2020).

642 13 Shen, J., McFarland, A. G., Young, V. B., Hayden, M. K. & Hartmann, E. M. Toward
643 Accurate and Robust Environmental Surveillance Using Metagenomics. *Frontiers in
644 Genetics* **12**, doi:10.3389/fgene.2021.600111 (2021).

645 14 Nayfach, S. & Pollard, K. S. Toward Accurate and Quantitative Comparative
646 Metagenomics. *Cell* **166**, 1103-1116, doi:10.1016/j.cell.2016.08.007 (2016).

647 15 Eisenhofer, R. *et al.* Contamination in Low Microbial Biomass Microbiome Studies: Issues
648 and Recommendations. *Trends in Microbiology* **27**, 105-117,
649 doi:<https://doi.org/10.1016/j.tim.2018.11.003> (2019).

650 16 McLaren, M. R., Willis, A. D. & Callahan, B. J. Consistent and correctable bias in
651 metagenomic sequencing experiments. *eLife* **8**, e46923, doi:10.7554/eLife.46923 (2019).

652 17 Martí, J. M. Recentrifuge: Robust comparative analysis and contamination removal for
653 metagenomics. *PLOS Computational Biology* **15**, e1006967,
654 doi:10.1371/journal.pcbi.1006967 (2019).

655 18 Davis, N. M., Proctor, D. M., Holmes, S. P., Relman, D. A. & Callahan, B. J. Simple
656 statistical identification and removal of contaminant sequences in marker-gene and
657 metagenomics data. *Microbiome* **6**, 226, doi:10.1186/s40168-018-0605-2 (2018).

658 19 Bowers, R. M. *et al.* Impact of library preparation protocols and template quantity on the
659 metagenomic reconstruction of a mock microbial community. *BMC Genomics* **16**, 856,
660 doi:10.1186/s12864-015-2063-6 (2015).

661 20 Mahnert, A. *et al.* Man-made microbial resistances in built environments. *Nature
662 Communications* **10**, 968, doi:10.1038/s41467-019-08864-0 (2019).

663 21 Minich, J. J. *et al.* KatharoSeq Enables High-Throughput Microbiome Analysis from Low-
664 Biomass Samples. *mSystems* **3**, e00218-00217, doi:10.1128/mSystems.00218-17 (2018).

665 22 Nocker, A., Cheung, C. Y. & Camper, A. K. Comparison of propidium monoazide with
666 ethidium monoazide for differentiation of live vs. dead bacteria by selective removal of

667 DNA from dead cells. *J Microbiol Methods* **67**, 310-320, doi:10.1016/j.mimet.2006.04.015
668 (2006).

669 23 Wang, Y. *et al.* Whole microbial community viability is not quantitatively reflected by
670 propidium monoazide sequencing approach. *Microbiome* **9**, 17, doi:10.1186/s40168-020-
671 00961-3 (2021).

672 24 Elizaquivel, P., Aznar, R. & Sanchez, G. Recent developments in the use of viability dyes
673 and quantitative PCR in the food microbiology field. *J Appl Microbiol* **116**, 1-13,
674 doi:10.1111/jam.12365 (2014).

675 25 Fittipaldi, M., Nocker, A. & Codony, F. Progress in understanding preferential detection
676 of live cells using viability dyes in combination with DNA amplification. *Journal of*
677 *Microbiological Methods* **91**, 276-289, doi:<https://doi.org/10.1016/j.mimet.2012.08.007>
678 (2012).

679 26 Morton, J. T. *et al.* Establishing microbial composition measurement standards with
680 reference frames. *Nat Commun* **10**, 2719, doi:10.1038/s41467-019-10656-5 (2019).

681 27 Liwinski, T., Leshem, A. & Elinav, E. Breakthroughs and Bottlenecks in Microbiome
682 Research. *Trends Mol Med* **27**, 298-301, doi:10.1016/j.molmed.2021.01.003 (2021).

683 28 Vandeputte, D. *et al.* Quantitative microbiome profiling links gut community variation to
684 microbial load. *Nature* **551**, 507-511, doi:10.1038/nature24460 (2017).

685 29 Hingamp, P. *et al.* Exploring nucleo-cytoplasmic large DNA viruses in Tara Oceans
686 microbial metagenomes. *ISME J* **7**, 1678-1695, doi:10.1038/ismej.2013.59 (2013).

687 30 Satinsky, B. M., Gifford, S. M., Crump, B. C. & Moran, M. A. in *Methods in Enzymology*
688 Vol. 531 (ed Edward F. DeLong) 237-250 (Academic Press, 2013).

689 31 Yamamoto, N., Hospodsky, D., Dannemiller, K. C., Nazaroff, W. W. & Peccia, J. Indoor
690 emissions as a primary source of airborne allergenic fungal particles in classrooms. *Environ*
691 *Sci Technol* **49**, 5098-5106, doi:10.1021/es506165z (2015).

692 32 Satinsky, B. M. *et al.* The Amazon continuum dataset: quantitative metagenomic and
693 metatranscriptomic inventories of the Amazon River plume, June 2010. *Microbiome* **2**, 17,
694 doi:10.1186/2049-2618-2-17 (2014).

695 33 Tkacz, A., Hortala, M. & Poole, P. S. Absolute quantitation of microbiota abundance in
696 environmental samples. *Microbiome* **6**, 110, doi:10.1186/s40168-018-0491-7 (2018).

697 34 Stammler, F. *et al.* Adjusting microbiome profiles for differences in microbial load by
698 spike-in bacteria. *Microbiome* **4**, 28, doi:10.1186/s40168-016-0175-0 (2016).

699 35 Rodriguez, R. L. & Konstantinidis, K. T. Nonpareil: a redundancy-based approach to
700 assess the level of coverage in metagenomic datasets. *Bioinformatics* **30**, 629-635,
701 doi:10.1093/bioinformatics/btt584 (2014).

702 36 Rodriguez, R. L. & Konstantinidis, K. T. Estimating coverage in metagenomic data sets
703 and why it matters. *ISME J* **8**, 2349-2351, doi:10.1038/ismej.2014.76 (2014).

704 37 Rodriguez-R, L. M., Gunturu, S., Tiedje, J. M., Cole, J. R. & Konstantinidis, K. T.
705 Nonpareil 3: Fast Estimation of Metagenomic Coverage and Sequence Diversity. *mSystems*
706 **3**, e00039-00018, doi:10.1128/mSystems.00039-18 (2018).

707 38 Leung, M. H., Wilkins, D., Li, E. K., Kong, F. K. & Lee, P. K. Indoor-air microbiome in
708 an urban subway network: diversity and dynamics. *Appl Environ Microbiol* **80**, 6760-6770,
709 doi:10.1128/AEM.02244-14 (2014).

710 39 Davis, A. *et al.* Improved yield and accuracy for DNA extraction in microbiome studies
711 **with variation in microbial biomass**. *BioTechniques* **66**, 285-289, doi:10.2144/btn-2019-
712 0016 (2019).

713 40 Hu, J. *et al.* Impacts of indoor surface finishes on bacterial viability. *Indoor Air* **29**, 551-
714 562, doi:10.1111/ina.12558 (2019).

715 41 Mo, Y. *et al.* Biogeographic patterns of abundant and rare bacterioplankton in three
716 subtropical bays resulting from selective and neutral processes. *ISME J* **12**, 2198-2210,
717 doi:10.1038/s41396-018-0153-6 (2018).

718 42 Nyirabuhoro, P. *et al.* Seasonal Variability of Conditionally Rare Taxa in the Water
719 Column Bacterioplankton Community of Subtropical Reservoirs in China. *Microbial*
720 *Ecology* **80**, 14-26, doi:10.1007/s00248-019-01458-9 (2020).

721 43 Ji, B. W. *et al.* Quantifying spatiotemporal variability and noise in absolute microbiota
722 abundances using replicate sampling. *Nature methods* **16**, 731-736, doi:10.1038/s41592-
723 019-0467-y (2019).

724 44 Barbau-Piednoir, E. *et al.* Evaluation of viability-qPCR detection system on viable and
725 dead *Salmonella* serovar Enteritidis. *J Microbiol Methods* **103**, 131-137,
726 doi:10.1016/j.mimet.2014.06.003 (2014).

727 45 Bengtsson-Palme, J. *et al.* Metaxa2 Database Builder: enabling taxonomic identification
728 from metagenomic or metabarcoding data using any genetic marker. *Bioinformatics* **34**,
729 4027-4033, doi:10.1093/bioinformatics/bty482 (2018).

730 46 Bengtsson-Palme, J., Thorell, K., Wurzbacher, C., Sjöling, Å. & Nilsson, R. H. Metaxa2
731 Diversity Tools: Easing microbial community analysis with Metaxa2. *Ecological*
732 *Informatics* **33**, 45-50, doi:<https://doi.org/10.1016/j.ecoinf.2016.04.004> (2016).

733 47 Bengtsson-Palme, J. *et al.* METAXA2: improved identification and taxonomic
734 classification of small and large subunit rRNA in metagenomic data. *Mol Ecol Resour* **15**,
735 1403-1414, doi:10.1111/1755-0998.12399 (2015).

736 48 Quast, C. *et al.* The SILVA ribosomal RNA gene database project: improved data
737 processing and web-based tools. *Nucleic Acids Res* **41**, D590-596,
738 doi:10.1093/nar/gks1219 (2013).

739 49 Yilmaz, P. *et al.* The SILVA and "All-species Living Tree Project (LTP)" taxonomic
740 frameworks. *Nucleic acids research* **42**, D643-D648, doi:10.1093/nar/gkt1209 (2014).

741 50 Beghini, F. *et al.* Integrating taxonomic, functional, and strain-level profiling of diverse
742 microbial communities with bioBakery 3. *eLife* **10**, e65088, doi:10.7554/eLife.65088
743 (2021).

744 51 Truong, D. T., Tett, A., Pasolli, E., Huttenhower, C. & Segata, N. Microbial strain-level
745 population structure and genetic diversity from metagenomes. *Genome Res* **27**, 626-638,
746 doi:10.1101/gr.216242.116 (2017).

747 52 Lax, S. *et al.* Bacterial colonization and succession in a newly opened hospital. *Sci Transl*
748 *Med* **9**, eaah6500, doi:10.1126/scitranslmed.aah6500 (2017).

749 53 O'Hara, N. B. *et al.* Metagenomic characterization of ambulances across the USA.
750 *Microbiome* **5**, 125, doi:10.1186/s40168-017-0339-6 (2017).

751 54 Chng, K. R. *et al.* Cartography of opportunistic pathogens and antibiotic resistance genes
752 in a tertiary hospital environment. *Nat Med* **26**, 941-951, doi:10.1038/s41591-020-0894-4
753 (2020).

754 55 Mishra, P. *et al.* Descriptive statistics and normality tests for statistical data. *Ann Card*
755 *Anaesth* **22**, 67-72, doi:10.4103/aca.ACA_157_18 (2019).

756 56 Blaustein, R. A. *et al.* Toothbrush microbiomes feature a meeting ground for human oral
757 and environmental microbiota. *Microbiome* **9**, 32, doi:10.1186/s40168-020-00983-x
758 (2021).

759 57 Glass, E. M. *et al.* MIxS-BE: a MIxS extension defining a minimum information standard
760 for sequence data from the built environment. *The ISME Journal* **8**, 1-3,
761 doi:10.1038/ismej.2013.176 (2014).

762 58 Adams, R. I., Bateman, A. C., Bik, H. M. & Meadow, J. F. Microbiota of the indoor
763 environment: a meta-analysis. *Microbiome* **3**, 49, doi:10.1186/s40168-015-0108-3 (2015).

764 59 El Khouli, R. H. *et al.* Relationship of temporal resolution to diagnostic performance for
765 dynamic contrast enhanced MRI of the breast. *J Magn Reson Imaging* **30**, 999-1004,
766 doi:10.1002/jmri.21947 (2009).

767 60 Escobar-Zepeda, A. *et al.* Analysis of sequencing strategies and tools for taxonomic
768 annotation: Defining standards for progressive metagenomics. *Sci Rep* **8**, 12034,
769 doi:10.1038/s41598-018-30515-5 (2018).

770 61 Marotz, C. A. *et al.* Improving saliva shotgun metagenomics by chemical host DNA
771 depletion. *Microbiome* **6**, 42, doi:10.1186/s40168-018-0426-3 (2018).

772 62 Nelson, M. T. *et al.* Human and Extracellular DNA Depletion for Metagenomic Analysis
773 of Complex Clinical Infection Samples Yields Optimized Viable Microbiome Profiles.
774 *Cell Rep* **26**, 2227-2240 e2225, doi:10.1016/j.celrep.2019.01.091 (2019).

775 63 Barlow, J. T., Bogatyrev, S. R. & Ismagilov, R. F. A quantitative sequencing framework
776 for absolute abundance measurements of mucosal and lumenal microbial communities. *Nat*
777 *Commun* **11**, 2590, doi:10.1038/s41467-020-16224-6 (2020).

778 64 Costea, P. I. *et al.* Towards standards for human fecal sample processing in metagenomic
779 studies. *Nat Biotechnol* **35**, 1069-1076, doi:10.1038/nbt.3960 (2017).

780 65 Pacholewicz, E. *et al.* Internal sample process control improves cultivation-independent
781 quantification of thermotolerant *Campylobacter*. *Food Microbiol* **78**, 53-61,
782 doi:10.1016/j.fm.2018.09.017 (2019).

783 66 Hardwick, S. A. *et al.* Synthetic microbe communities provide internal reference standards
784 for metagenome sequencing and analysis. *Nat Commun* **9**, 3096, doi:10.1038/s41467-018-
785 05555-0 (2018).

786 67 Turlousse, D. M. *et al.* Synthetic spike-in standards for high-throughput 16S rRNA gene
787 amplicon sequencing. *Nucleic Acids Res* **45**, e23, doi:10.1093/nar/gkw984 (2017).

788 68 Shen, J. & Hartmann, E. M. Collect, extract, pool, and cultivate surface swab samples from
789 built environments. *Protocol Exchange*, doi:10.21203/rs.3.pex-1656/v1 (2021).

790 69 Gomez-Silvan, C. *et al.* A comparison of methods used to unveil the genetic and metabolic
791 pool in the built environment. *Microbiome* **6**, 71, doi:10.1186/s40168-018-0453-0 (2018).

792 70 Shen, J. & Hartmann, E. M. Snap freezing of environmental microbial samples in liquid.
793 *Protocol Exchange*, doi:10.21203/rs.3.pex-1657/v1 (2021).

794 71 Wang, X., Howe, S., Deng, F. & Zhao, J. Current Applications of Absolute Bacterial
795 Quantification in Microbiome Studies and Decision-Making Regarding Different
796 Biological Questions. *Microorganisms* **9**, 1797 (2021).

797 72 Fahimipour, A. K. *et al.* Daylight exposure modulates bacterial communities associated
798 with household dust. *Microbiome* **6**, 175, doi:10.1186/s40168-018-0559-4 (2018).

799 73 Shen, J., Rose, L. J., Perry-Dow, K. A. & Hartmann, E. M. An optimized protocol for
800 propidium monoazide treatment. *Protocol Exchange*, doi:10.21203/rs.3.pex-1659/v1
801 (2021).

802 74 Leung, M. H. Y. *et al.* Characterization of the public transit air microbiome and resistome
803 reveals geographical specificity. *Microbiome* **9**, 112, doi:10.1186/s40168-021-01044-7
804 (2021).

805 75 Shen, J. & Hartmann, E. M. DNA extraction protocol for low-biomass environmental
806 samples: adapted from the Lucigen MasterPure Complete DNA and RNA Purification Kit
807 manual. *Protocol Exchange*, doi:10.21203/rs.3.pex-1658/v1 (2021).

808 76 Singer, V. L., Jones, L. J., Yue, S. T. & Haugland, R. P. Characterization of PicoGreen
809 reagent and development of a fluorescence-based solution assay for double-stranded DNA
810 quantitation. *Analytical biochemistry* **249**, 228-238, doi:10.1006/abio.1997.2177 (1997).

811 77 Klindworth, A. *et al.* Evaluation of general 16S ribosomal RNA gene PCR primers for
812 classical and next-generation sequencing-based diversity studies. *Nucleic Acids Res* **41**, e1,
813 doi:10.1093/nar/gks808 (2013).

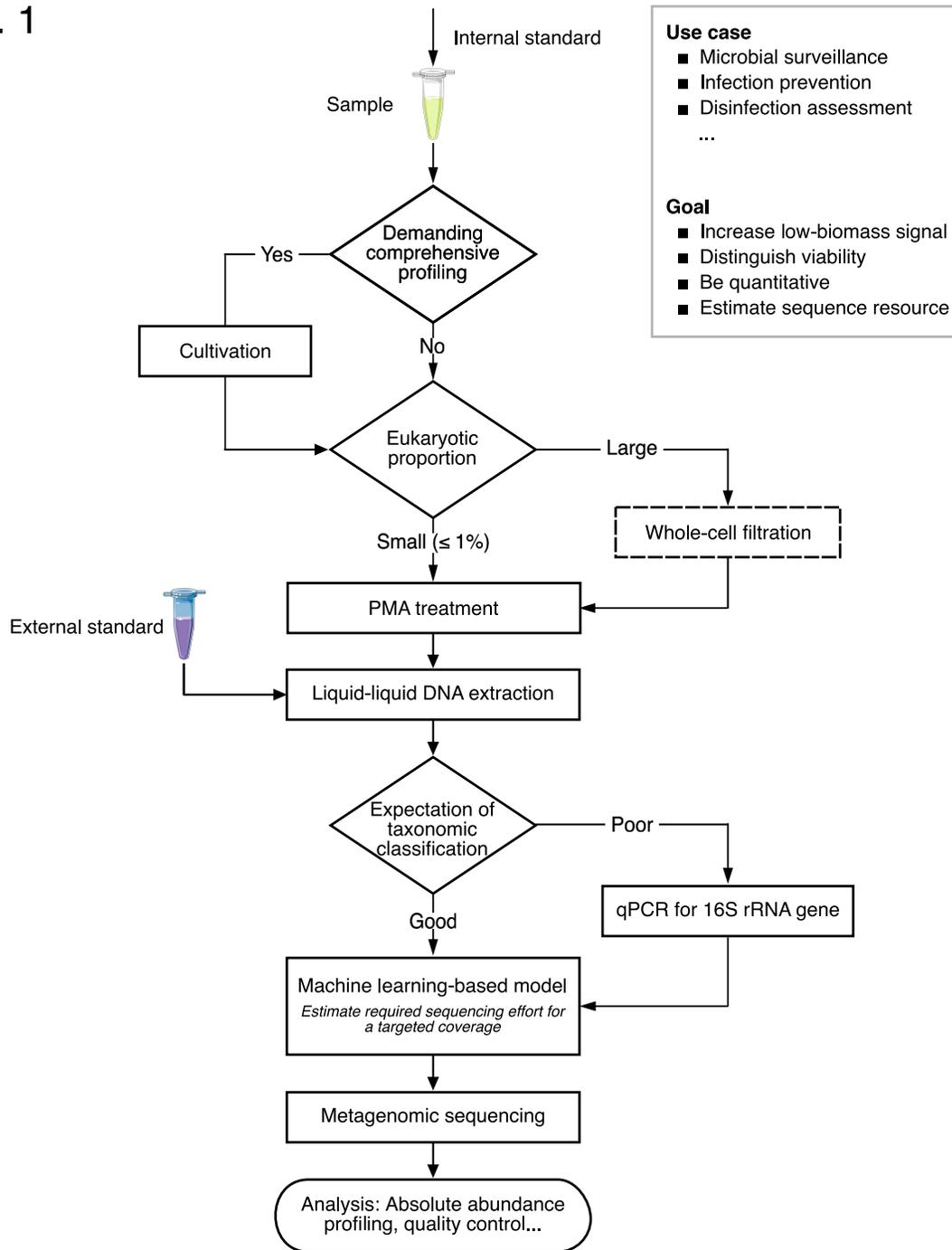
814 78 Segata, N. *et al.* Metagenomic microbial community profiling using unique clade-specific
815 marker genes. *Nature Methods* **9**, 811-814, doi:10.1038/nmeth.2066 (2012).

816 79 Team, R. C. R: A language and environment for statistical computing. (2013).

817 80 Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*. (Springer-Verlag New York,
818 2016).

819

Fig. 1



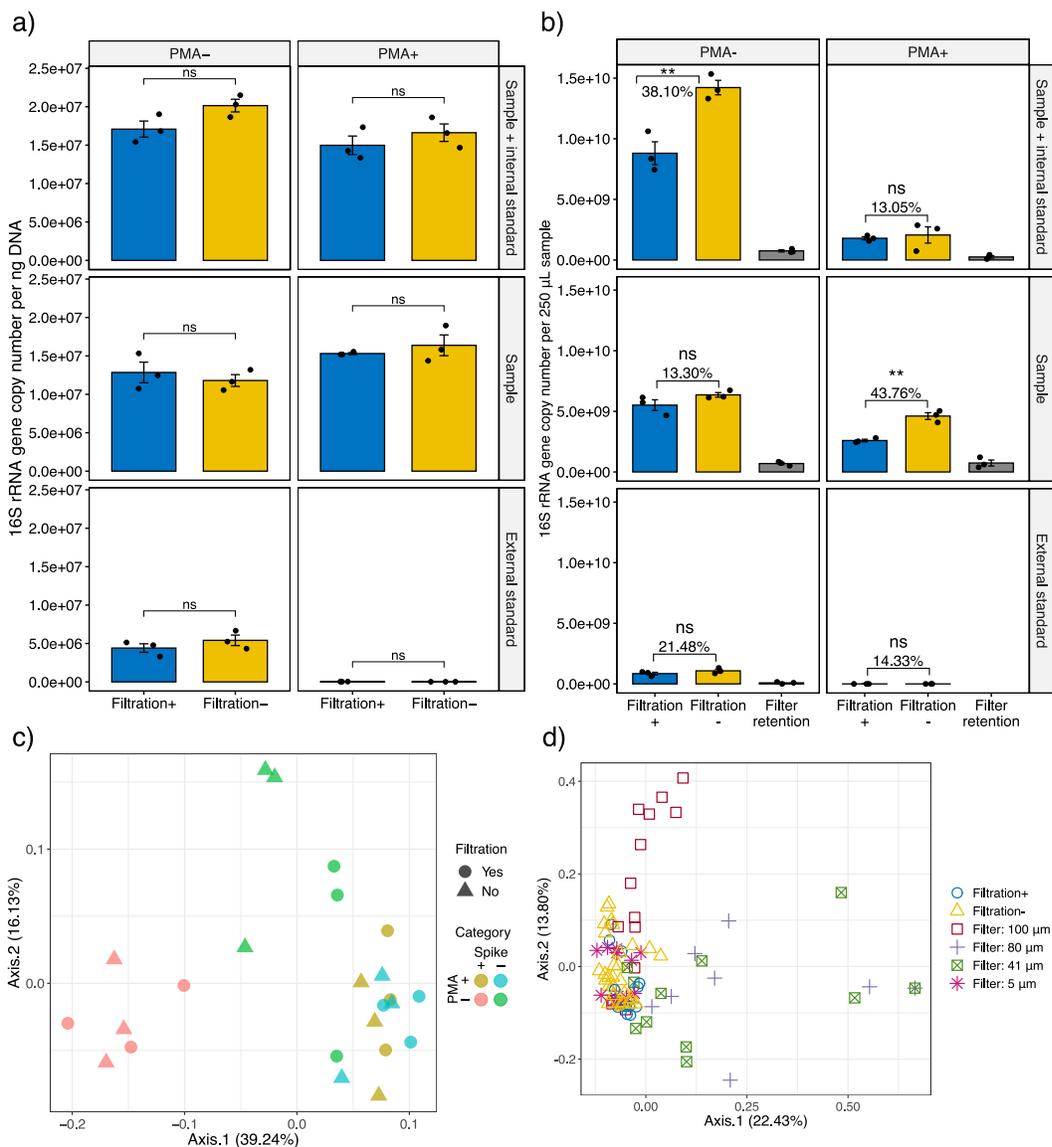
Treatment in dashed box needs further investigation.

820

821 **Figure 1: A workflow for metagenomics-based environmental surveillance that is appropriate for low-biomass**
 822 **samples, distinguishes viability, is quantitative, and estimates sequence resources.**

823

Fig. 2

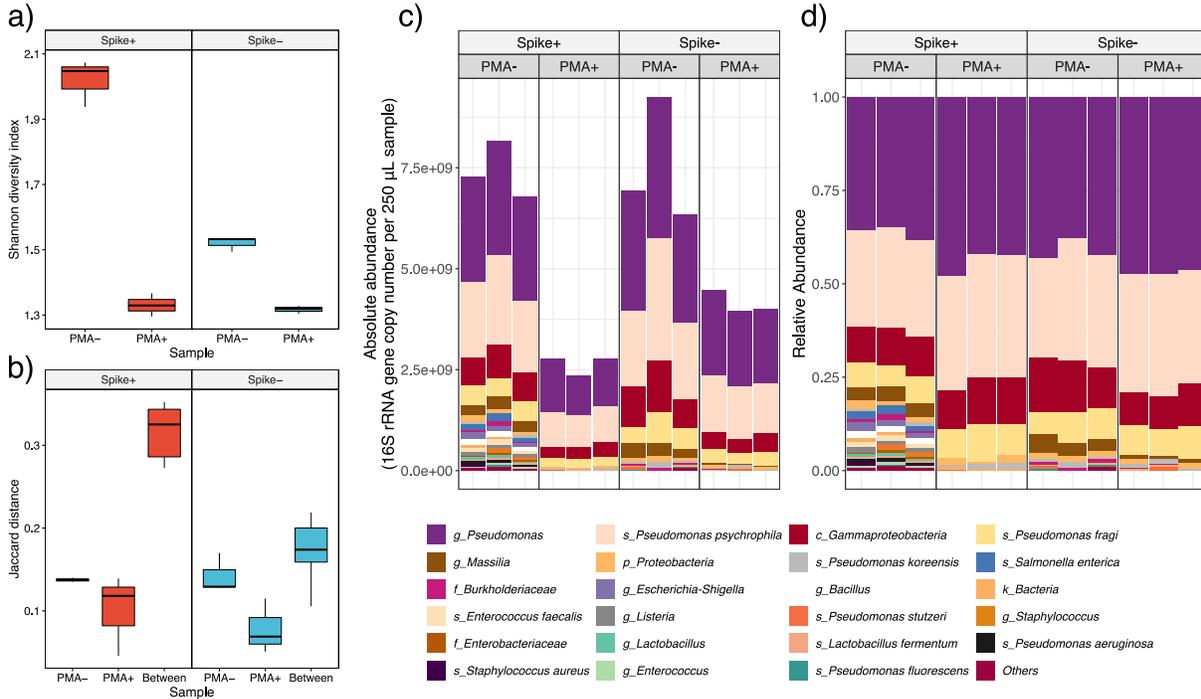


824

825 **Figure 2: Effects of sequential filtration on hospital-associated surface samples.** a) Bacteria proportion was not
 826 significantly increased after filtration according to paired t tests. b) Biomass of samples with and without filtration as
 827 well as retained by filters according to 16S rRNA gene copy number. In a) and b), error bars represent the mean
 828 standard error of triplicates. Filter retention includes all biomass captured by 100, 80, 41, and 5 µm filters. Ns and **
 829 are significance codes, representing $p > 0.05$ and $0.001 < p \leq 0.01$, respectively. A linear scale was used for both a)
 830 and b) because for a), a linear scale is more conservative than a log scale when no significant difference was concluded;
 831 for b), linear-scale biomass loss is more informative for metagenomic sequencing. c) Principle coordinate analysis
 832 using Jaccard distance metric among samples with and without filtration. d) Principle coordinate analysis based on
 833 Jaccard distance metric revealed that bacterial profiles retained on 5 µm filters clustered together with liquid samples,
 834 while those on 100, 80, and 41 µm filters were away from the major group.

835

Fig. 3

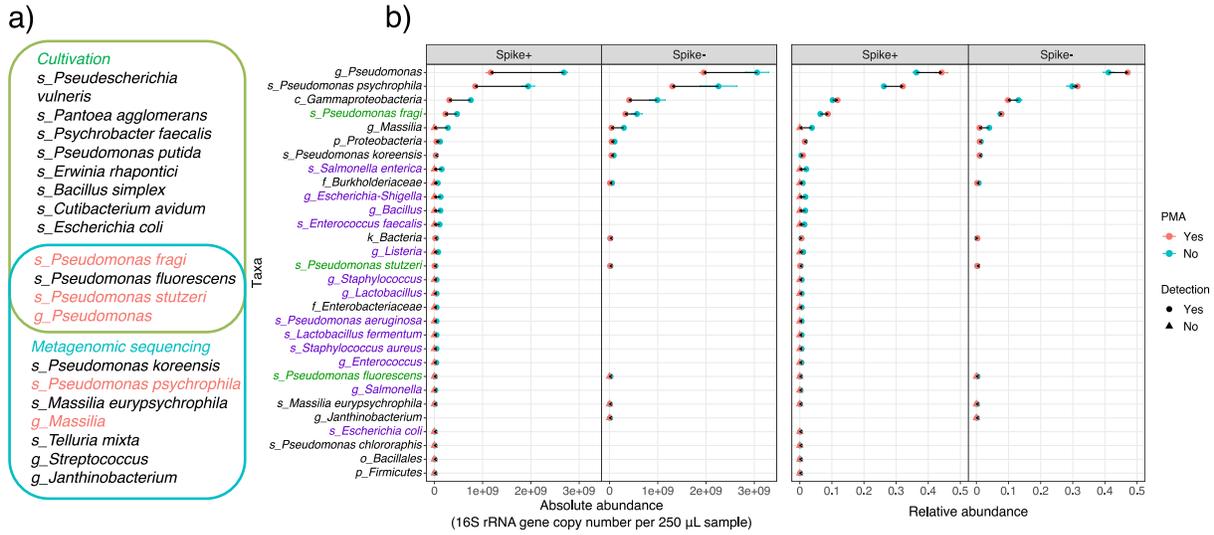


836

837 **Figure 3: Effects of PMA treatment on hospital-associated surface samples.** a) PMA treated samples had lower α
 838 diversity based on Shannon index. b) Inter distances between paired samples with and without PMA treatment were
 839 larger than intra distances within each sample group (based on Jaccard metric). Comparisons of profiling the bacterial
 840 composition by c) absolute abundance and c) relative abundance.

841

Fig. 4

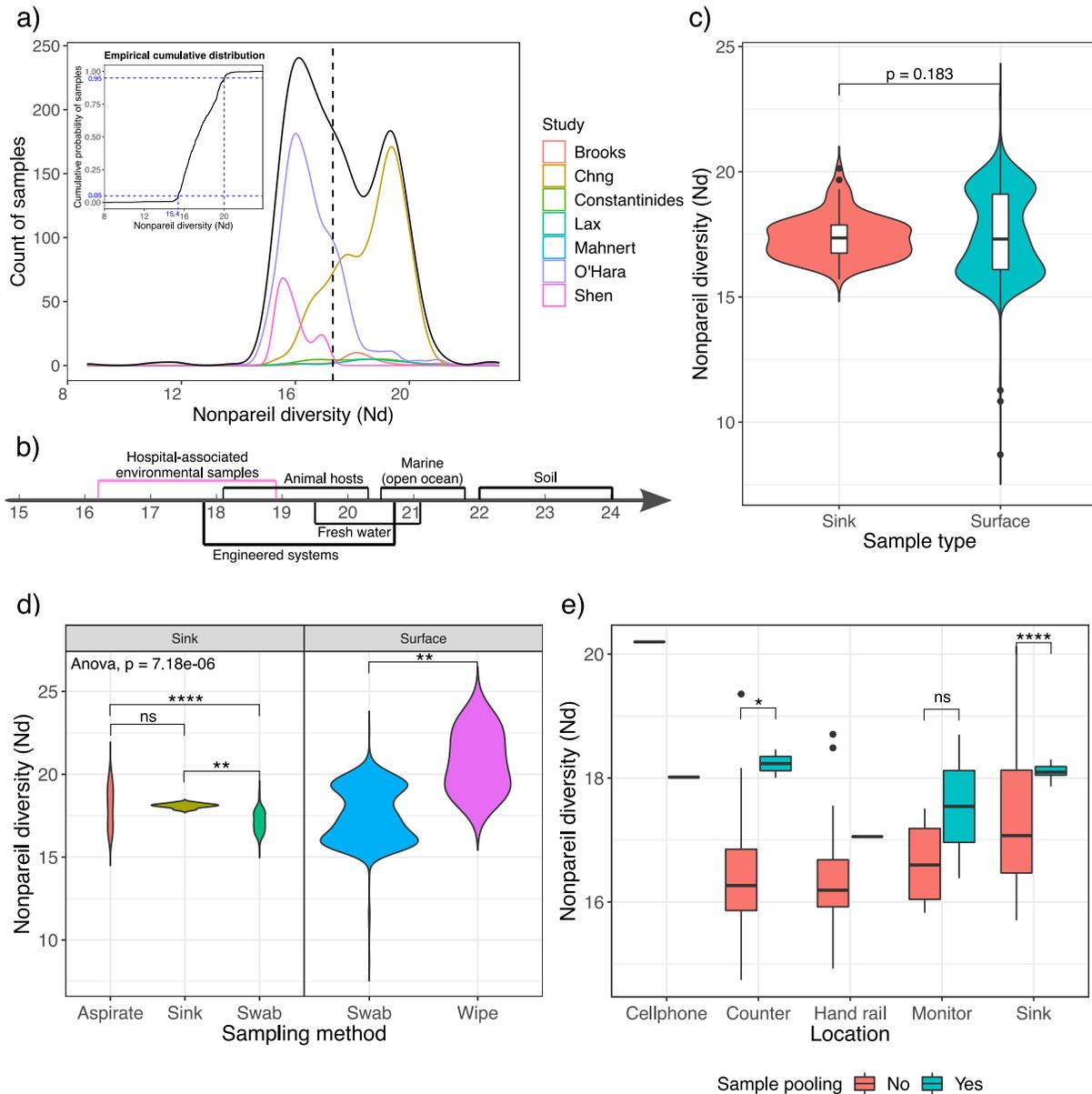


842

843 **Figure 4: Performance of cultivation and PMA-MetaSeq in viability distinction of hospital-associated surface**
 844 **samples.** a) Venn diagram showing the detected taxa by cultivation and MetaSeq. Taxa detected by PMA-MetaSeq
 845 are color coded in red. b) The abundance change of all taxa detected under the framework of absolute abundance and
 846 relative abundance. Taxa in the theoretical composition of the internal standard and recovered in cultivation are color
 847 coded in purple and green, respectively. The Y axis follows a descending order of the average abundance across
 848 samples. Error bars represent the mean standard error of triplicates.

849

Fig. 5

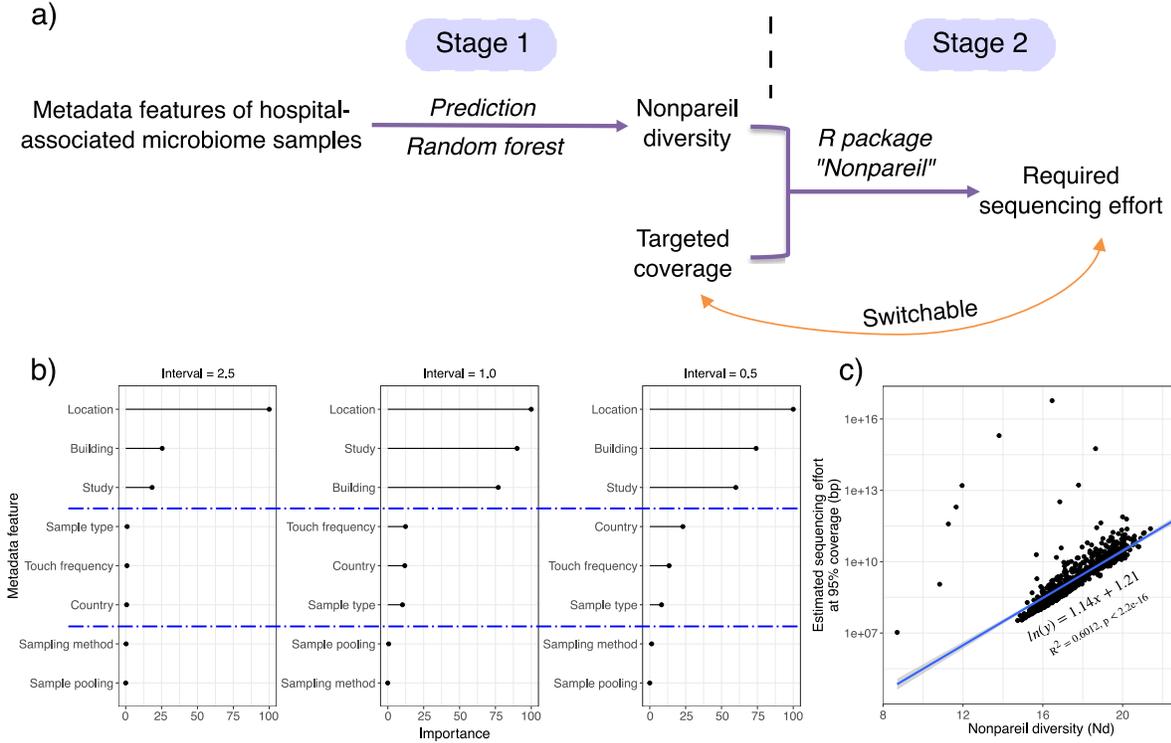


850

851 **Figure 5: Relationships between Nonpareil diversity and metadata features for hospital-associated surface**
 852 **samples.** a) Overall distribution of Nonpareil diversity (black) and distributions for individual studies. b) Interquartile
 853 range of Nonpareil diversity for microbiome samples from different environments. This study is color coded in orchid.
 854 Effects of c) sample type, d) sampling method, e) and sample pooling on Nonpareil diversity. Significance codes are:
 855 $p > 0.05$ (ns), $0.01 < p \leq 0.05$ (*), $0.001 < p \leq 0.01$ (**), $p \leq 0.001$ (****).

856

Fig. 6



857

858 **Figure 6: Required sequencing effort can be predicted by accessible sample features and targeted coverage.** a)
 859 Workflow of making the prediction. b) Variable importance rankings by random forest. c) The natural log of estimated
 860 sequencing effort at 95% coverage is linearly correlated with Nonpareil diversity.

Supplementary Files

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

- [supt2samplingcondition.xlsx](#)
- [supt1predictiondataset.xlsx](#)
- [supplementaryfigurescombined.pdf](#)