

1 **Clinical practices underlie COVID-19 patient respiratory microbiome composition and**
2 **its interactions with the host**

3

4 Verónica Lloréns-Rico^{1,2}, Ann C. Gregory^{1,2}, Johan Van Weyenbergh³, Sander Jansen⁴,
5 Tina Van Buyten⁴, Junbin Qian^{5,6}, Marcos Braz³, Soraya Maria Menezes³, Pierre Van
6 Mol^{5,6,7}, Lore Vanderbeke⁸, Christophe Doods^{7,9}, Jan Gunst¹⁰, Greet Hermans¹⁰, Philippe
7 Meersseman¹¹, CONTAGIOUS collaborators, Els Wauters^{7,9}, Johan Neyts⁴, Diether
8 Lambrechts^{5,6}, Joost Wauters^{11,12}, Jeroen Raes^{1,2,12,13}

9

10 ¹ Laboratory of Molecular Bacteriology, Department of Microbiology and Immunology,
11 Rega Institute, KU Leuven, Belgium

12 ² Center for Microbiology, VIB, Leuven, Belgium

13 ³ Laboratory for Clinical and Evolutionary Virology, Department of Microbiology and
14 Immunology, Rega Institute, KU Leuven, Belgium

15 ⁴ Laboratory of Virology and Chemotherapy, Department of Microbiology, Immunology
16 and Transplantation, Rega Institute, KU Leuven, Belgium

17 ⁵ Laboratory of Translational Genetics, Department of Human Genetics, KU Leuven,
18 Belgium

19 ⁶ VIB Center for Cancer Biology, VIB, Leuven, Belgium

20 ⁷ Department of Pneumology, University Hospitals Leuven, Belgium

21 ⁸ Laboratory of Clinical Bacteriology and Mycology, Department of Microbiology,
22 Immunology and Transplantation, KU Leuven, Belgium

23 ⁹ Laboratory of Respiratory Diseases and Thoracic Surgery (BREATHE), Department of
24 Chronic Diseases and Metabolism, KU Leuven, Belgium

25 ¹⁰ Laboratory of Intensive Care Medicine, Department of Cellular and Molecular
26 Medicine, KU Leuven, Belgium

27 ¹¹ Laboratory for Clinical Infectious and Inflammatory Disorders, Department of
28 Microbiology, Immunology and Transplantation, KU Leuven, Belgium

29 ¹² These authors contributed equally

30 ¹³ Corresponding author: jeroen.raes@kuleuven.vib.be

31

32 **Keywords:** COVID-19, SARS-CoV-2, respiratory microbiome, single-cell RNA-sequencing,
33 host-microbiome interactions

34

35 **Abstract**

36

37 Understanding the pathology of COVID-19 is a global research priority. Early evidence
38 suggests that the respiratory microbiome may be playing a role in disease progression,
39 yet current studies report contradictory results. Here, we examine potential
40 confounders in COVID-19 respiratory microbiome studies by analyzing the upper (n=58)
41 and lower (n=35) respiratory tract microbiome in well-phenotyped COVID-19 patients
42 and controls combining microbiome sequencing, viral load determination, and
43 immunoprofiling. We found that time in the intensive care unit and the type of oxygen
44 support, both of which are associated to additional treatments such as antibiotic usage,
45 explained the most variation within the upper respiratory tract microbiome, while SARS-
46 CoV-2 viral load had a reduced impact. Specifically, mechanical ventilation was linked to
47 altered community structure, lower species- and higher strain-level diversity, and
48 significant shifts in oral taxa previously associated with COVID-19. Single-cell
49 transcriptomic analysis of the lower respiratory tract of mechanically ventilated COVID-
50 19 patients identified specific oral bacteria, different to those observed in controls.
51 These oral taxa were found physically associated with proinflammatory immune cells,
52 which showed higher levels of inflammatory markers. Overall, our findings suggest
53 confounders are driving contradictory results in current COVID-19 microbiome studies
54 and careful attention needs to be paid to ICU stay and type of oxygen support, as
55 bacteria favored in these conditions may contribute to the inflammatory phenotypes
56 observed in severe COVID-19 patients.

57

58 **Introduction**

59

60 COVID-19, a novel coronavirus disease classified as a pandemic by the World Health
61 Organization, has caused over 150 million reported cases and 3 million deaths
62 worldwide to date. Infection by its causative agent, the novel coronavirus SARS-CoV-2,
63 results in a wide range of clinical manifestations: it is estimated that around 80% of
64 infected individuals are asymptomatic or present only mild respiratory and/or
65 gastrointestinal symptoms, while the remaining 20% develop acute respiratory distress
66 syndrome requiring hospitalization and oxygen support and, of those, 25% of cases

67 necessitate critical care. Despite a concerted global research effort, many questions
68 remain about the full spectrum of the disease severity. Independent studies from
69 different countries, however, agree that age and sex are the major risk factors for
70 disease severity and patient death¹⁻³, as well as type 2 diabetes and obesity^{4,5}. Other
71 potential risk factors for critical condition and death are viral load of the patient upon
72 hospital admission⁶⁻⁸ and the specific immune response to infection, with manifestation
73 of an abnormal immune response in critical patients characterized by dysregulated
74 levels of pro-inflammatory cytokines and chemokines, which some studies have
75 associated with organ failure^{9,10}.

76

77 Despite its close interplay with the immune system and its known associations with host
78 health, little is known about the role of the respiratory microbiota in modulating COVID-
79 19 disease severity, or its potential as a prognostic marker¹¹. Previous studies exploring
80 other pulmonary disorders have shown that lung microbiota members may exacerbate
81 symptoms and contribute to their severity¹², potentially through direct crosstalk with
82 the immune system and/or due to bacteremia and secondary infections¹³. First studies
83 of the respiratory microbiome in COVID-19 have revealed elevated levels of
84 opportunistic pathogenic bacteria¹⁴⁻¹⁶. However, reports on bacterial diversity are
85 contradictory. While some studies report a low microbial diversity in COVID-19
86 patients^{14,17} that rebounds following recovery¹⁵, others show an increased diversity in
87 the COVID-19 associated microbiota¹⁶. These conflicting results could be due to
88 differences in sampling location (upper or lower respiratory tract), severity of the
89 patients, disease stage, treatment or other confounders. While these early findings
90 already suggest that the lung microbiome could be exacerbating or mitigating COVID-19
91 progression, exact mechanisms are yet to be elucidated. Therefore, an urgent need
92 exists for studies identifying and tackling confounders in order to discern true signals
93 from noise.

94

95 To identify potential associations between COVID-19 severity and evolution and the
96 upper and lower respiratory tract microbiota, we used nasopharyngeal swabs and
97 bronchoalveolar lavage (BAL) samples, respectively. For the upper respiratory tract, we
98 longitudinally profiled the nasopharyngeal microbiome of 58 COVID-19 patients during

99 intensive care unit (ICU) treatment and after discharge to a classical hospital ward
100 following clinical improvement, in conjunction with viral load determination and
101 nCounter immune profiling. For the lower respiratory tract, we profiled microbial reads
102 in cross-sectional single-cell RNA-seq data¹⁸ from of bronchoalveolar lavage (BAL)
103 samples of 22 COVID-19 patients and 13 pneumonitis controls with negative COVID-19
104 qRT-PCR, obtained from the same hospital. The integration of these data enabled us to
105 (1) identify potential confounders of COVID-19 microbiome associations, (2) explore
106 how microbial diversity evolves throughout hospitalization, (3) study microbe-host cell
107 interaction and (4) substantiate a link between the respiratory microbiome and SARS-
108 CoV-2 viral load as well as COVID-19 disease severity. Altogether, our results suggest the
109 existence of associations between the microbiota and specific immune cells in the
110 context of COVID-19 disease. These interactions may be driven by mechanical
111 ventilation and its associated clinical practices, and therefore could potentially influence
112 COVID-19 disease progression and resolution.

113

114 **Results**

115

116 **The upper respiratory microbiota of COVID-19 patients**

117

118 We longitudinally profiled the upper respiratory microbiota of 58 patients diagnosed
119 with COVID-19 based on a positive qRT-PCR test or a negative test with high clinical
120 suspicion based on symptomatology and a chest CT-scan showing typical round glass
121 opacities. All these patients were admitted and treated at UZ Leuven hospital. Patient
122 demographics for this cohort are shown in Table 1.

123

124 In total, 112 nasopharyngeal swabs from these patients were processed (Figure 1a): the
125 V4 region of the 16S rRNA gene amplified on extracted DNA using 515F and 806R
126 primers, and sequenced on an Illumina MiSeq platform (see Methods). From the same
127 swabs, RNA was extracted to determine SARS-CoV-2 viral loads and to estimate immune
128 cell populations of the host and expression of immune-related genes using nCounter
129 (Methods). Of the 112 samples processed and sequenced, 101 yielded over 10,000
130 amplicon reads that could be assigned to bacteria at the genus level (Figure 1b;

131 Methods). The microbiome of the entire cohort was dominated by the gram-positive
132 genera *Staphylococcus* and *Corynebacterium*, typical from the nasal cavity and
133 nasopharynx¹⁹.

134

135 **Bacterial alpha diversity is associated with ICU stay length, SARS-CoV-2 viral load and** 136 **calprotectin levels**

137

138 First, we determined genus-level alpha-diversity for the 101 samples with more than
139 10,000 genus-level assigned reads, using the Shannon Diversity index (SDI; see Methods;
140 Supplementary Table 1). We observed that the SDI was significantly different across
141 sampling moments (Kruskal-Wallis test, p-value = 0.009; Supplementary Figure 1a), with
142 significant differences between swabs procured upon patient ICU admission and later
143 timepoints, suggesting an effect of disease progression and/or treatment (for instance
144 due to antibiotics administered throughout ICU stay). We explored these differences
145 further, and observed that SDI correlated negatively with the number of days spent in
146 ICU at the moment of sampling, with longer ICU stays leading to a lower diversity ($\rho=-$
147 0.53, p-value= $1.9 \cdot 10^{-8}$).

148

149 To evaluate the association of other clinical or disease-related variables with upper
150 respiratory tract microbiome diversity, we used a generalized linear mixed model
151 framework: we performed an exhaustive screening of all possible models containing up
152 to 8 different explanatory variables, using an automated model selection algorithm (see
153 Methods). The variables used to regress the SDI comprise the patient ID, modeled as a
154 random effect; disease-related variables, such as the time in ICU, SARS-CoV-2 viral load
155 or the use of mechanical ventilation; and other variables known to affect the
156 microbiome, such as the administration of antibiotics (specifically
157 meropenem/piperacillin-tazobactam and ceftriaxone) or the levels of inflammatory
158 markers (calprotectin, C-reactive protein). The antibiotics meropenem and
159 piperacillin/tazobactam were grouped as a single variable in all subsequent analyses as
160 they were administered under the same clinical guidelines. The best performing model
161 (AICc=121.79; p-value= $4.06 \cdot 10^{-8}$) included the patient modeled as a random effect and
162 confirmed a negative association between the time spent in ICU and diversity.

163 Additionally, this model showed a negative effect of SARS-CoV-2 viral load and a positive
164 association of calprotectin levels with the SDI (Figure 1c,d; Supplementary Figure 1b-d).

165

166 We leveraged all the models generated in the screening to calculate weighted
167 importance scores for all the fixed effects tested (Methods; Supplementary Figure 1e).

168 These scores showed that the three variables incorporated in the best model (time in
169 ICU, SARS-CoV-2 viral load and calprotectin) held the highest relative importance,
170 followed by CRP levels and mechanical ventilation. Treatment with antibiotics
171 ceftriaxone an meropenem or piperacillin-tazobactam had the lowest importance
172 scores, and no significant differences in SDI were found between samples obtained
173 before and after the administration of meropenem/piperacillin-tazobactam
174 (Supplementary Figure 1f).

175 Altogether, our data suggest that respiratory microbiome diversity is linked to the length
176 of ICU stay, SARS-CoV-2 viral load and calprotectin levels. While no significant effects
177 were found for the most widely used antibiotics in this cohort, we cannot rule out that
178 antibiotic administration or other clinical practices are causing the decrease of SDI over
179 time.

180

181

182 **Respiratory microbiome composition variation is linked to respiratory support and**
183 **associated clinical practices**

184

185 We next explored potential associations between the upper respiratory genus-level
186 microbiota composition and the extensive metadata collected in the study. In total, 72
187 covariates related to patient anthropometrics, medication and clinical variables
188 measured in the hospital, as well as SARS-CoV-2 viral load, host cytokine expression and
189 estimated immune cell populations measured in the swabs were tested (Supplementary
190 Table 2). Individually, 20 of these covariates showed a significant correlation to
191 microbiota composition in a univariate analysis (dbRDA, p -value<0.05; FDR<0.05; Figure
192 2a). These significant covariates were related to disease and measures of its severity,
193 such as the clinical evaluation of the patient, the total length of the ICU stay, the number
194 of days in ICU at the time of sampling, or the type of oxygen support needed by the

195 patient. Despite showing an association to the overall diversity, SARS-CoV-2 viral load
196 detected in the swabs was not significantly associated to microbiome composition
197 variation (Supplementary Table 2). Neither ongoing antibiotic usage (i.e., administration
198 of any type of antibiotic) nor number of ongoing antibiotics administered were
199 significant, but the administration of specific antibiotics meropenem/piperacillin-
200 tazobactam (previous or ongoing treatment) and ceftriaxone (ongoing administration
201 only) showed significant associations with microbiome composition (Supplementary
202 Table 2, Figure 2a).

203

204 Of the 20 significant covariates, only 2 accounted for 48.7% non-redundant variation in
205 this dataset in a multivariate analysis (dbRDA; p-value=0.001), with the rest holding
206 redundant information. These were the patient ID, included due to the longitudinal
207 sampling of patients, and confirming that intra-individual variation over time is smaller
208 than patient inter-individual variation²⁰, and the type of oxygen support received at the
209 time of sampling (Figure 2a,b). Notably, the type of oxygen support discriminated
210 samples based on ventilation type, with non-invasive ventilation samples (groups 1, 2
211 and 3) separating from samples from intubated patients (groups 4 to 7; PERMANOVA
212 test, $R^2=0.0642$, p-value=0.001). Because of this separation, we also evaluated whether
213 previous mechanical ventilation (regardless of the specific group) had a significant
214 impact on the microbiome composition, showing even a larger effect size than when
215 considering only the ongoing mechanical ventilation (PERMANOVA test, $R^2=0.0965$, p-
216 value=0.001), suggesting that this invasive procedure may have an effect that is
217 prolonged in time.

218

219 Mechanical ventilation is inherently associated to additional clinical practices, such as
220 administration of broad-spectrum antibiotics and decontamination procedures
221 (including chlorhexidine washes) to prevent/treat ventilator-associated pneumonia.
222 Hence, we explored whether antibiotic usage could explain the significant relationship
223 between microbiome composition and oxygen support type. We found that from the
224 specific antibiotics associated to microbiome composition, ceftriaxone was
225 predominantly administered in patients on non-invasive oxygen support (Chi-square, p-
226 value=0.001), whilst meropenem or piperacillin-tazobactam were preferentially given to

227 patents on mechanical ventilation (Chi-square test, p-value=0.002; Supplementary
228 Figure 2a). This association is not casual and responds to current treatment guidelines
229 at UZ Leuven: ceftriaxone is administered to patients upon admission and for 3-7 days
230 to prevent potential bacterial co-infections. In our cohort, 80% of the patients received
231 ceftriaxone at the beginning of their stay (Supplementary Figure 2b). Patients with
232 longer ICU stays and requiring higher levels of oxygen support will be considered to have
233 hospital-acquired/ventilator-associated pneumonia (HAP/VAP) and receive meropenem
234 or piperacillin/tazobactam (Supplementary Figure 2b). Therefore, the observed
235 correlation between oxygen support types and these antibiotics can be explained by
236 disease severity and length of ICU stay.

237

238 We therefore explored whether we could observe an effect of oxygen support type
239 alone, deconfounding for the patient ID and the two significant antibiotic covariates
240 using partial dbRDA to extract the effect size of oxygen support alone. The
241 deconfounded model exhibited a significant association to overall microbiome
242 composition (partial dbRDA; $R^2=0.058$, p-value=0.042) suggesting that although
243 antibiotic administration may explain part of the variation in microbiome composition
244 observed, there may an independent effect of the oxygen support type. Nevertheless,
245 the effect of other practices concomitant to mechanical ventilation, such as oral
246 decontamination with chlorhexidine washes, could not be disentangled as these
247 treatments were always performed together.

248

249 To determine if oxygen support or associated practices also impacted the microbiome
250 at finer taxonomic resolution, we revisited alpha-diversity at species- and strain-level.
251 We defined species as 97% identity 16S OTUs and strains per species as the clustered
252 16S sequences within each OTU. Our analyses revealed both species- and strain-level
253 diversity change with ventilation, even with non-invasive ventilation (e.g. BIPAP, CPAP).
254 Across all samples we observed high species- and low strain-level diversity pre-
255 ventilation, which reversed following any form of ventilation (Figure 2c; Wilcoxon test;
256 p-values<0.05, with the exception of type 7), with the exception of ventilation with
257 inhaled nitric oxide. Further, species- and strain-level diversity showed a strong inverse
258 correlation (Figure 2d; Pearson's correlation, $R^2 = -0.92$, p-value = 0.0035).

259

260 Given the observed effect of mechanical ventilation on the overall microbiome
261 composition, we evaluated which specific taxa were differentially abundant between
262 samples from intubated and non-intubated patients. In total, 28 genera were more
263 abundant in samples from mechanically ventilated patients, while 1 genus was more
264 abundant in non-invasively ventilated patients (p -value<0.05; FDR<0.05; Figure 2e,
265 Supplementary Figure 3a; Supplementary Table 3). When controlling for the effect of
266 the antibiotics ceftriaxone and meropenem/piperacillin, 20 genera were significantly
267 different between both groups of samples (Supplementary Figure 3b, Supplementary
268 Table 3). Some of these taxa are common oral microbiome commensals or opportunistic
269 pathogens that had been repeatedly reported as more abundant in COVID-19 patients
270 than in controls, such as *Prevotella*, *Fusobacterium*, *Porphyromonas* or *Lactobacillus*^{14–}
271 ¹⁶. Here, we reported higher abundance of these genera in mechanically ventilated
272 COVID-19 patients as compared to non-mechanically ventilated COVID-19 patients. This
273 points at mechanical ventilation (and associated practices such as oral decontamination)
274 as a potential confounder of previous COVID-19 studies. Additionally, we found other
275 taxa not previously reported in previous COVID-19 microbiome studies, such as
276 *Mycoplasma* or *Megasphaera* (Figure 2e, Supplementary Figure 2), but previously
277 associated to risk of ventilator-associated pneumonia²¹.

278

279 By extracting the amplicon sequence variants (ASVs) corresponding to these
280 differentially abundant genera (see Methods), some of these taxa could be narrowed
281 down to the species level, confirming their origin as typically oral bacteria: for instance,
282 *Prevotella* species included *P. oris*, *P. salivae*, *P. denticola*, *P. buccalis* and *P. oralis*.
283 Within the *Mycoplasma* genus, ASVs were assigned to *Mycoplasma salivarium* among
284 other species, an oral bacterium which has been previously associated to the incidence
285 of ventilator-associated pneumonia²¹. When controlling for ventilation type, no taxa
286 were found associated to SARS-CoV-2 viral loads (Supplementary Table 3). These results
287 show that further research with larger cohorts and controlling for the relevant
288 confounders highlighted here, such as ventilation type, antibiotic usage or length of stay
289 in ICU, will be needed to study the specific effect of the viral infection.

290

291 **Single-cell RNA-seq of bronchoalveolar fluid identifies oral commensals and**
292 **opportunistic pathogens in the lower respiratory tract**

293

294 Next, we explored what the functional consequences of (disease and/or treatment-
295 driven) lung microbiome disturbances could be. To do so, we screened host single-cell
296 RNA-seq data generated on BAL samples of 35 patients¹⁸ using a computational pipeline
297 to identify microbial reads (see Methods). All patients in this cross-sectional cohort
298 showed clinical symptoms of pneumonia, 22 of them being diagnosed with COVID-19.
299 The other 13 patients with non-COVID-19 pneumonia were hereafter referred to as
300 controls (Table 1). Out of the 35 patients, 21 were admitted to ICU (20 COVID-19 patients
301 and 1 control) and 14 were hospitalized in ward at the moment of sampling (2 COVID-
302 19 patients and 12 controls; Table 1 and Supplementary Figure 4). Microbial read
303 screening of these samples revealed an average of 7,295.3 microbial reads per sample
304 (ranging from 0 to 74,226 reads, with only a single sample yielding zero microbial reads;
305 Supplementary Figure 4).

306

307 Among the top taxa encountered in these patients, we found similarities with the data
308 obtained in nasopharyngeal swabs. The top 15 species detected include *Mycoplasma*
309 *salivarium* as the dominating taxon in 5 COVID-19 patients in ICU, as well as different
310 *Prevotella* members. Non-COVID-19 pneumonia patients in ward (i.e. non-mechanically
311 ventilated) harbored different microbes: 2 patients had a microbiome dominated by
312 *Porphyromonas gingivalis*, while a single patient had a microbiome dominated by the
313 fungus *Pneumocystis jirovecii*, a known pathogen causing *Pneumocystis* pneumonia
314 (PCP)²².

315

316 Supplementary table 4 shows associations between organism abundances and specific
317 patient metadata: disease, hospital stay and ventilation type. Multiple links with COVID-
318 19 diagnosis were identified (Wilcoxon test, (noncorrected) p-value<0.05; see Methods)
319 but due to the low sample number, none was significant after multiple-test correction.
320 Additionally, as hospital stay (ICU or ward), type of oxygen support (invasive or non-
321 invasive ventilation) and disease (COVID-19 or controls) were highly correlated (Chi-
322 squared test, p-value < 0.0001 for all three pairwise correlations), the effect of these

323 three variables could not be disentangled. Therefore, although this data may validate
324 our findings from the upper respiratory tract microbiome, due to the small cohort size
325 and the existence of multiple confounders, these association results should be
326 confirmed in larger studies.

327

328 **Bacteria in the lower respiratory tract associate to host cells from the innate immune** 329 **system in COVID-19 patients**

330

331 Next, we took advantage of the single-cell barcoding and questioned whether the
332 microbial reads that we identified were found in association with host cells (for instance
333 infecting or internalized), or contrarily, had unique barcodes suggesting a free-living
334 state. In total, 29,886 unique barcodes were identified that matched a total of 46,151
335 microbial UMIs. The distribution of UMIs per barcode was asymmetrical, ranging from 1
336 to 201 and with 88% of the barcodes having a single UMI. Additionally, 26,572 barcodes
337 (89%) were associated to a single microbial species, the rest being associated to 2
338 species (8.8%) or more (2.2%).

339

340 Out of the total 29,886 microbial barcodes, only 2,108 were also assigned to host cells,
341 suggesting that the bulk of bacteria found in BAL samples exist as free-living organisms
342 or in bacterial biofilms. Although microscopic evaluation would be needed to validate
343 this hypothesis, bacterial biofilms have been previously documented in bronchoalveolar
344 lavages²³, and the enrichment for host cells in these samples via centrifugation¹⁸ may
345 have also indirectly enriched these specimens for biofilm and/or host-associated
346 microbes. However, for the fraction of bacteria associated to host cells, the distribution
347 across disease types was not random. We found that while 2.3% of the non-COVID-19
348 patient cells were associated to bacterial cells, almost the double (4%) could be
349 observed in COVID-19 patients (Figure 3a; Chi-squared test; p-value < $2.2 \cdot 10^{-16}$).
350 However, because COVID-19 diagnosis is highly correlated with mechanical ventilation
351 in this cohort, this effect could be due to higher intubation rates in COVID-19 patients
352 and possibly, a higher incidence of VAP. Within COVID-19 patients, we also evaluated
353 the overlap between bacteria-associated host cells and cells with detected SARS-CoV-2
354 reads¹⁸ (Supplementary Table 5). Out of 1,033 host cells associated with bacteria in

355 COVID-19 patients and 343 cells with detected SARS-CoV-2 reads, only one cell was
356 positive for both viral and bacterial reads. A binomial test for independence of virus and
357 bacteria detection in the same host cell, showed that the observed co-occurrence in one
358 cell only was highly unlikely ($p\text{-value}=5.7\cdot 10^{-4}$), therefore suggesting mutual exclusion of
359 microbiome members and viruses in the same host immune cells. However, it must be
360 noted that lack of detection does not necessarily imply lack of association of
361 bacteria/virus to host cells, especially with experimental methods such as single-cell
362 RNA-seq, designed for profiling host cells and not optimized for detection of these
363 entities. Therefore, further studies with larger sample sizes are required to validate the
364 co-exclusion hypothesis.

365

366 We also explored whether host-associated bacterial reads would preferentially be linked
367 with specific cell types, taking into account the varying frequencies of cell types in
368 COVID-19 patients and controls (see Methods). Such a preferential association would
369 suggest that these observations are biologically relevant and not an artifact of the single-
370 cell sample and library preparation. Among control patients, cell types were similarly
371 distributed in both groups (i.e. with and without bacteria), with only a preferential
372 association of microbial cells with neutrophils ($p\text{-value} = 3.61\cdot 10^{-12}$; Figure 3b;
373 Supplementary Figure 4). However, in COVID-19 patients, three cell types were
374 significantly associated with bacteria: neutrophils ($p\text{-value} < 2.2\cdot 10^{-16}$), monocytes ($p\text{-}$
375 $\text{value} = 4.82\cdot 10^{-5}$) and monocyte-derived macrophages ($p\text{-value} < 2.2\cdot 10^{-16}$; Figure 3b;
376 Supplementary Figure 5). We also found that different bacteria associate with distinct
377 host cells. For instance, in COVID-19 patients, bacteria from the *Mycoplasma* genus
378 preferentially associated to monocyte-derived macrophages ($p\text{-value} = 2.28\cdot 10^{-7}$), while
379 *Rothia* ($p\text{-value} = 8.21\cdot 10^{-4}$), *Enterobacter* ($p\text{-value} = 2.59\cdot 10^{-5}$), or *Klebsiella* ($p\text{-value} =$
380 $3.12\cdot 10^{-9}$) are enriched in monocytes (Figure 3c).

381

382 Last, we investigated whether the associations of bacteria to host cells are linked to host
383 cell expression. To do so, we assessed whether expression based cell subtype
384 classification¹⁸ for neutrophils, monocytes and macrophages showed non-random
385 associations with bacteria across all samples in this cohort. Among the neutrophils, a
386 subtype of inflammatory neutrophils characterized by expression of the calgranulin

387 S100A12 was enriched in bacteria-associated cells (p-value $7.18 \cdot 10^{-6}$; Figure 3d,e). This
388 subset of cells was also found to be enriched in SARS-CoV-2 nucleocapsid gene reads¹⁸,
389 suggesting that the same cell type responsible for defense against the virus would be
390 responding to potentially invasive bacteria in the lung. This subgroup is characterized by
391 the expression of the calprotectin subunits S100A8 and S100A9. It is known that
392 S100A8/A9 heterodimer secretion is increased in infection-induced inflammation and
393 has some antibacterial effects mediated by secretion of pro-inflammatory cytokines,
394 release of reactive oxygen species and recruitment of other inflammatory cells, as well
395 as chelation of Zn^{2+} necessary for bacterial enzymatic activity²⁴. These mechanisms are
396 mediated by binding of the S100A8/A9 dimer to TLR4 receptors to trigger the release of
397 pro-inflammatory cytokines such as IL-6 and TNF- α , and thus may contribute to sustain
398 or exacerbate inflammation²⁵. Therefore, the association with bacteria could, at least in
399 part, explain the inflammatory phenotype of this neutrophil subset. To further examine
400 this hypothesis, we explored differential gene expression between bacteria-associated
401 and non-associated S100A12^{hi} neutrophils (Supplementary Table 6). Because
402 association of these cells with SARS-CoV-2 and with bacteria was mutually exclusive, we
403 also compared these changes with the ones triggered by the virus in neutrophils²⁶.
404 Within this subset, neutrophils with co-occurring bacteria showed significantly higher
405 expression (Bonferroni-corrected p-value < 0.05) of pro-inflammatory genes, including
406 the cytokine IL1B and some of its target genes (PTSG2), the transcription factors FOS
407 and JUN, and several genes involved in degranulation (S100A9, FOLR3, HSPA1A,
408 HSP90AA1, FCGR3B), (Supplementary Table 6). Among these, FOLR3, a gene encoding
409 for a folate receptor, is found in neutrophil secretory granules and has antibacterial
410 functions, by binding folates and thus depriving bacteria of these essential
411 metabolites²⁷. This response differed to that of virus-engulfing neutrophils in that IFN
412 response genes are not distinctively upregulated by bacteria.

413

414 Regarding myeloid cells, both inflammatory IL1B^{hi} monocytes (p-value = $2 \cdot 10^{-16}$) as well
415 as a mixed group of CCL2-expressing macrophages (p-value = $5.38 \cdot 10^{-10}$) are enriched in
416 bacteria-associated cells (Figure 3f). These inflammatory monocytes are believed to
417 have an important role in the aberrant immune response occurring in severe COVID-19
418 patients. In this case, further gene expression patterns were detected, specific for

419 bacteria-associated cells: for CCL2^{hi} macrophages, cells with co-occurring bacteria
420 expressed higher levels of MHC genes of type I and II, suggesting a more active role of
421 these cells in antigen presentation (Bonferroni-corrected p-value < 0.05; Figure 3f;
422 Supplementary Table 6). A similar increase was also observed in monocytes, yet not
423 significant (Supplementary Table 6), possibly due to the lower monocyte abundances in
424 this dataset. Additionally, bacteria-associated macrophages express significantly higher
425 levels of the calprotectin subunits S100A8/A9, similarly to neutrophils, as well as pro-
426 inflammatory chemokines (such as CCL4, CXCL10 and CXCL1).

427

428 Altogether, our results suggest that the bacteria detected in these cell subsets via
429 scRNA-seq analyses may be contributing to the inflammatory response observed in the
430 host.

431

432 **Discussion**

433

434 Since the beginning of the COVID-19 pandemic, a massive global effort by the scientific
435 community was undertaken to understand physiopathology of SARS-CoV-2 infection
436 and risk factors affecting disease outcome. In this study, we explored the respiratory
437 microbiota as a potential risk factor for disease severity, and we evaluated the upper
438 and lower respiratory tract microbiota in COVID-19 patients throughout hospitalization.
439 We linked this data to viral load measurements and immunoprofiling results from
440 nCounter and single-cell RNA sequencing data. To assess robustness of previously
441 reported signals, we investigated the effect of potential confounders based on a broad
442 panel of patient metadata variables.

443

444 We found that in the upper respiratory tract, SARS-CoV-2 viral load has a mild negative
445 association with bacterial biodiversity. A larger effect of severity indicators such as
446 calprotectin levels or length of ICU stay was observed, with diversity decreasing
447 throughout the length of the ICU period, a pattern reminiscent of that seen in other
448 pulmonary conditions^{28,29}. The effect of ICU length-of-stay may be mediated by
449 treatment options such as the administration of broad-spectrum antibiotics and/or
450 patient intubation and mechanical ventilation. Antibiotic usage might also explain why

451 calprotectin levels correlate with alpha diversity: such antibiotics would decrease overall
452 microbial diversity, including (some of) the taxa that could be linked to inflammation.
453 The observed effects of these clinical practices on microbiome alpha diversity could
454 potentially explain why previous studies on the microbiota of COVID-19 patients show
455 conflicting results regarding diversity: some studies reported lower diversity in sputum
456 or throat swab samples of COVID-19 patients^{14,15,17} while others focusing on the lower
457 respiratory microbiome using bronchoalveolar fluid samples, showed higher bacterial
458 diversity in COVID-19 patients than in controls¹⁶. To further complicate matters, it
459 cannot be excluded that sampling site or processing could also be potential confounders
460 in these studies and/or reflect the different pathologies in the different areas of the
461 respiratory tract.

462

463 We further found that between-patient microbiome variation (as measured by genus-
464 level microbial beta-diversity) was also influenced by different severity indicators such
465 as the clinical status of the patient, or more importantly the type of oxygen support
466 received, with mechanically ventilated patients harboring a different microbiota than
467 non-intubated patients. This effect could not be fully explained by neither general
468 antibiotic administration, nor the usage of specific antibiotics such as ceftriaxone,
469 meropenem or piperacillin-tazobactam, suggesting an independent effect of mechanical
470 ventilation. Such an independent effect has previously been suggested in small
471 cohorts^{28,30,31}, but it needs to be validated in larger studies. However, other associated
472 practices such as decontamination procedures could still be responsible for the
473 observed associations. The impact of oxygen support was also reflected at the species-
474 and strain-levels, with intubation causing a significant decrease and increase,
475 respectively, in diversity. We hypothesize that the introduction of forced oxygen may
476 drive the fast extinction of certain microbial species enabling the diversification of
477 existing or newly colonizing species into new strains. Combined, these results suggest
478 that non-invasive ventilation (e.g. BIPAP, CPAP) can have microbial effects indicating
479 that any form of ventilation may be a tipping point for microbial community differences.

480

481 Importantly, several of the taxa reported to change between intubated and non-
482 intubated patients were reported to be linked to diagnosis in previous COVID-19

483 microbiome studies^{14–16}. In our study, no taxa were specifically linked to SARS-CoV-2
484 viral load after controlling for mechanical ventilation. This result suggests the possibility
485 that mechanical ventilation and its associated clinical practices are confounding
486 previous results. Indeed, one study comparing COVID-19 patients with patients
487 diagnosed of community-acquired pneumonia found no differences in respiratory
488 microbiome composition between both groups of patients, but both groups did differ
489 from healthy controls³². Together, these results indicate that patient intubation or even
490 non-invasive ventilation, as well as their associated medical practices, are to be
491 considered as important confounders when studying the upper respiratory microbiome,
492 and we strongly suggest future COVID-19 microbiome studies should foresee and
493 include strategies to account for this covariate. As an example, a recent study found a
494 single ASV corresponding to the genus *Rothia* that was specific for SARS-CoV-2 patients
495 after controlling for ICU-related confounders by comparing with a previous study of the
496 microbiome in ICU patients³³. Additionally, these findings on potential drivers of
497 microbiome variation are not exclusive to COVID-19 disease: the effect of intubation on
498 the respiratory microbiome and its influence on the incidence of ventilator-associated
499 pneumonia have been previously studied^{28,30,31}.

500

501 To better understand the potential functional consequences of these procedures and
502 linked microbial shifts, we also profiled the microbiome of the lower respiratory tract
503 using single-cell data obtained from a cross-sectional cohort of patients derived from
504 the same hospital. Our results show that single-cell RNA-seq, despite not being
505 optimized for microbial detection and profiling, can identify bacteria alone or in
506 association with specific human cells. Unfortunately, the low numbers of microbial reads
507 obtained in this small cohort, together with the fact that ICU stay, COVID-19 diagnosis
508 and intubation are highly correlated in this set of patients, only allow for a first
509 exploratory analysis of the results, requiring validation in further datasets. In this cohort,
510 we identified different oral commensals and opportunistic pathogens previously linked
511 to COVID-19 patients in both groups of samples, thus pointing again at a potential
512 ventilation-linked origin. More interestingly, we identified a subset of bacteria
513 associated with host cells, more specifically with neutrophils, monocytes and
514 macrophages. This enrichment shows that these bacteria are likely not random

515 contaminants, from which an even distribution across cell types (i.e. considering cell
516 type abundances) would be expected. The identity of these host cells suggests that
517 bacteria could have been phagocytosed by these innate immune system cells, rather than
518 be attached to the host cell surface. To the best of our knowledge, this is the first study
519 linking interacting host cells and lung microbiome via high-throughput single-cell RNA-
520 seq.

521

522 We find that host cells associated with bacteria, most of which are of oral origin, exhibit
523 pro-inflammatory phenotypes as well as higher levels of MHC for antigen presentation.
524 In this single-cell cohort it was observed that critical COVID-19 patients are characterized
525 by an impaired monocyte to macrophage differentiation, resulting in an excess of pro-
526 inflammatory monocytes, as well as by prolonged neutrophil inflammation²⁶. Given that
527 only these cell types are enriched in bacteria, we hypothesize that the respiratory (or
528 ventilation-linked) microbiome may be playing a role in exacerbating COVID-19
529 progression in the lower respiratory tract. We verified that this response could likely be
530 driven by bacteria and not SARS-CoV-2, which is also detected mostly in these cell types,
531 as there is almost no overlap in detection of both virus and bacteria in the same cells.
532 However, it must be noted that lack of detection does not completely rule out presence
533 of virus and bacteria within these cells. Therefore, further research is required in order
534 to confirm a causative role of the microbiota in this immune impairment characteristic
535 of critical disease, and to reveal the specific mechanisms involved.

536

537 The presence of oral taxa in the lung microbiota is not unique of disease conditions. It is
538 known that microaspiration, or the aspiration of aerosol droplets originated in the oral
539 cavity, occurs in healthy individuals and can serve as a route for lung colonization of oral
540 commensals³⁴. Such an increase of oral bacteria in the lower respiratory tract could be
541 facilitated when critically ill patients –including but not limited to COVID-19– get
542 intubated. Indeed, oral bacteria have been linked to ventilator-associated
543 pneumonia^{35,36}. It is yet to be elucidated whether COVID-19 physiopathology favors lung
544 colonization by oral bacteria or if, in contrast, a lung microbiome previously colonized
545 by oral microbes could also contribute to the disease. What is known is that an increase

546 of oral bacteria in the lower respiratory tract can result in an increased inflammatory
547 phenotype, even in healthy subjects³⁷

548

549 **Conclusion**

550

551 Overall, this study provides a systematic analysis of potential confounders in COVID-19
552 microbiome studies. We identified that ICU hospitalization and type of oxygen support,
553 which may be at least partially explained by clinical practices such as antibiotic usage,
554 had large impacts on the upper respiratory tract microbiome and have the potential to
555 confound clinical microbiome studies. Among the different types of oxygen support we
556 reported the largest shifts in microbial community structure between intubated and
557 non-intubated patients. We found that oral taxa were strongly enriched in the upper
558 respiratory tract of mechanically ventilated COVID-19 patients, and specific taxa were
559 also found in the lower respiratory tract of COVID-19 patients. Further, in the lower
560 respiratory tract, microbes were strongly associated with specific pro-inflammatory
561 immune cells. This information contributes to a collective body of literature on the
562 pathology of COVID-19 and suggests that careful attention be paid to ICU stay and type
563 of oxygen support and associated clinical practices such as antibiotic usage or oral
564 decontamination procedures when evaluating the role of the lung microbiome on
565 COVID-19 disease progression.

566

567 **Methods**

568

569 **Study design and patient cohorts**

570

571 All experimental protocols and data analyses were approved by the Ethics Commission
572 from the UZ Leuven Hospital, under the COntAGlouS observational clinical trial (study
573 number S63381). The study design is compliant with all relevant ethical regulations,
574 including the Declaration of Helsinki and in the GDPR. All participants gave their
575 informed consent to participate in the study.

576

577 A total of 58 patients from the COntAGlouS observational trial were included as our
578 upper respiratory tract cohort. All patients were admitted to the UZ Leuven hospital
579 with a diagnostic of COVID-19. The disease was diagnosed based on a) a positive qRT-
580 PCR test, performed on admission or previously on other hospitals, when patients were
581 transferred from other medical facilities; or b) a chest CT-scan showing alveolar damage
582 and clinical symptoms of the disease. All patients included in the study were admitted
583 to ICU for a variable amount of time. Nasopharyngeal swabs were taken from these
584 patients at different timepoints throughout ICU stay and after ICU discharge, during
585 recovery in ward. A total of 112 swabs were processed for upper respiratory microbiome
586 characterization (Figure 1a).

587

588 To extend our findings from the upper respiratory tract, we also profiled the lower
589 respiratory tract microbiota in a different cohort¹⁸ of 35 patients belonging to the same
590 observational trial and also recruited at UZ Leuven hospital. This cross-sectional cohort
591 is composed by 22 COVID-19 patients and 13 pneumonitis controls with negative qRT-
592 PCR for SARS-CoV-2, with varying disease severity. Previous data from single-cell RNA-
593 sequencing had been collected for this cohort¹⁸. We reanalyzed this single-cell dataset
594 to profile the lower respiratory tract microbiota in these patients.

595

596

597 **RNA/DNA extraction and sequencing**

598

599 Nucleic acid extraction from the swab samples was performed with AllPrep
600 DNA/RNA/miRNA Universal kit (QIAGEN, catnr. 80224). Briefly, swabs from the
601 potentially infectious samples were inactivated by adding 600 μ L RLT-plus lysis buffer. To
602 increase bacterial cell lysis efficiency, glass beads and DX reagent (Pathogen Lysis Tubes,
603 QIAGEN, catnr. 19091) were added to the lysis buffer, and samples were disrupted in a
604 FastPrep-24TM instrument with the following program: 1-minute beating at 6.5m/sec, 1-
605 minute incubation at 4°C, 1-minute beating at 6.5m/sec, 1-minute incubation at 4°C.
606 After lysis, the remaining extraction steps followed the recommended protocol from the
607 manufacturer. DNA was eluted in 50 μ L EB buffer. Amplification of the V4 region of the
608 16S gene was done with primers 515F and 806R, using single multiplex identifiers and
609 adaptors as previously described³⁸. RNA was eluted in 30 μ L of nuclease-free water and
610 used for SARS-CoV-2 viral load determination in the swabs as well as to measure
611 inflammatory markers and cytokines and to estimate host cell populations via marker
612 gene expression using nCounter. In brief, raw nCounter data were processed using
613 nSolver 4.0 software (Nanostring), sequentially correcting three factors for each
614 individual sample: technical variation between samples (using spiked positive control
615 RNA), background correction (using spiked negative control RNA) and RNA content
616 variation (using 15 housekeeping genes). We have previously validated nCounter digital
617 transcriptomics for simultaneous quantification of host immune and viral transcripts³⁹,
618 including respiratory viruses in nasopharyngeal aspirates, even with low RNA yield⁴⁰⁻⁴².
619
620 DNA sequencing was performed on an Illumina MiSeq instrument, generating paired-
621 end reads of 250 base pairs.

622
623 For quality control, reads were demultiplexed with LotuS v1.565⁴³ and processed
624 following the DADA2 microbiome pipeline using the R packages DADA2⁴⁴ and
625 phyloseq⁴⁵. Briefly, reads were filtered and trimmed using the parameters truncQ=11,
626 truncLen=c(130,200), and trimLeft=c(30, 30) and then denoised. After removing
627 chimeras, amplicon sequence variants (ASVs) table was constructed and taxonomy was
628 assigned using the Ribosomal Database Project (RDP) classifier implemented in DADA2
629 (RDP trainset 16/release 11.5). The abundance table was then corrected for copy
630 number, rarefied to even sequencing depth, and decontaminated. For decontamination,

631 we used the prevalence-based contaminant identification method in the R package
632 decontam⁴⁶.

633

634 **16S statistical analysis**

635

636 All the 16S data analyses were performed using R v3.6.0 and the packages vegan
637 (v2.5.7)⁴⁷, phyloseq (v1.34.0)⁴⁵, CoDaSeq (v0.99.6)⁴⁸, DESeq2 (v1.30.1)⁴⁹, Biostrings
638 (v2.58.0)⁵⁰, rstatix (v0.7.0)⁵¹, glmulti (v1.0.8)⁵², sjPlot (v2.8.7)⁵³, and DECIPHER
639 (v2.18.1)⁵⁴.

640

641 To analyze the 16S amplicon data, technical replicates were pooled and counts from
642 technical replicates were added. For all the analyses using genus-level agglomerated
643 data, only samples containing more than 10,000 reads assigned at the genus level were
644 used (101 samples in total). Alpha-diversity was analyzed using Shannon's Diversity
645 Index. Comparison of the alpha diversity values across different groups was performed
646 using Kruskal-Wallis tests for comparisons across multiple groups. When applicable,
647 pairwise comparisons were performed using Dunn post-hoc tests. To establish the
648 potential associations of alpha diversity with different metadata variables, we selected
649 8 variables related to COVID-19 disease and/or known to affect microbiome
650 composition and diversity: patient ID, days spent in ICU, SARS-CoV-2 viral load, antibiotic
651 usage for ceftriaxone and meropenem/piperacillin-tazobactam, previous mechanical
652 ventilation, calprotectin gene expression and CRP levels. Meropenem and piperacillin-
653 tazobactam were merged as a single antibiotic as their administration is indicated under
654 the same clinical guidelines.

655

656 We used the R package glmulti to perform an exhaustive evaluation of the 256 models
657 including all possible combinations of the selected variables. All models generated were
658 generalized linear models or generalized linear mixed models (when including the
659 patient ID as a random effect), using a Gaussian family with a logarithmic link. Model
660 ranking and selection was performed based on the lowest small-sample-corrected
661 Akaike Information Criterion (AICc), and model significance was assessed comparing
662 with a null model (including the intercept only) using ANOVA test. Final variable

663 importance was calculated as a weighted average of the models in which each of the
664 variables appeared, with weights corresponding to the model ranks, defined by their
665 AICc values. This was also implemented as part of the glmulti package. The final model
666 plots were generated with the sjPlot package. Intra-patient differences in alpha diversity
667 between timepoints before and after administration of antibiotics or mechanical
668 ventilation were determined with Wilcoxon signed-rank tests.

669

670 Beta diversity analyses were performed using distance-based redundancy analyses
671 (dbRDA), using Aitchison distances. Prior to CLR data transformation, we filtered the
672 data using the CoDaSeq.filter function, to keep samples with more than 10,000 reads
673 and taxa with a relative abundance above 0.1% in any sample, as well as a prevalence of
674 at least 10% in the cohort. To replace zeros, we first calculated the minimum (non-zero)
675 relative abundance of each taxon across all samples. Then, for samples with zero counts
676 for a given taxon, the minimum relative abundance of the specific taxon was multiplied
677 by the total counts of such samples and this value was used to impute the zeros. dbRDA
678 analyses were performed using the capscale function from vegan, first in univariate
679 analyses with 72 metadata variables (Supplementary Table 2). Model p-values were
680 corrected using Benjamini-Hochberg's (BH) multiple-testing correction, to select 20
681 variables with BH-adjusted p-values < 0.05. These 20 variables were included in a
682 multivariate model, and non-redundant contribution to variation was calculated using
683 forward stepwise variable selection via the ordiR2step function from vegan. To
684 deconfound the effect of antibiotics and patient ID for oxygen support type, partial
685 dbRDA was used, including both antibiotics and patient ID as blocking variables.
686 Metadata variables containing dates, as well as non-informative metadata were
687 excluded. Non-informative metadata variables were defined as those containing a single
688 non-NA value or, for categorical variables, those being unevenly distributed (with >90%
689 of the samples belonging to the same category, for instance an antibiotic administered
690 only in two different samples). Additionally, from pairs of highly collinear variables
691 (correlation higher than 0.9), only one variable was kept.

692

693 Differential taxa abundance analyses were performed using DESeq2's likelihood ratio
694 tests and controlling for potential confounders when indicated, including them in the
695 null model.

696

697 To explore species-level and strain-level diversity, 16S sequences were first clustered
698 into 97% nucleotide diversity operational taxonomic units (OTUs) using the R packages
699 Biostrings and DECIPHER. These OTUs were used to represent the species-level. The
700 number of unique 16S sequences clustered within each OTU were used to represent the
701 number of detectable strains per species. To calculate strain-level diversity per sample,
702 the number of strains of 5 detected OTU species were randomly selected and averaged.
703 This was repeated 1,000x and the average of the all 1,000 subsamplings was used as the
704 final strain-level diversity value for each sample, as previously described⁵⁵. To account
705 for uneven sampling assessing diversity differences based on different parameters, we
706 randomly selected and averaged the species- and strain-level diversity of 5 samples per
707 parameter. This was repeated 100x and the subsamplings were used to assess the
708 significant differences between species- and strain-level diversity across the
709 parameters. The average was of all 100 subsamplings was used to as the input for a
710 Pearson's correlation between species- and strain-level diversity.

711

712 All statistical tests were two-sided unless otherwise specified, and when multiple tests
713 were applied to the different features (e.g. the differential taxa abundances across two
714 conditions) p-values were corrected for multiple testing using Benjamini-Hochberg's
715 method.

716

717 **Identification of microbial reads in BAL scRNA-seq data**

718 BAL scRNA-seq raw fastq data, as well as cell type and subtype assignments for all
719 individual cells, were obtained from a previous publication from within the COntAGlous
720 consortium. Experimental procedures on BAL samples as well as detailed host single-cell
721 gene expression analyses are detailed in the original publication¹⁸.

722

723 5' single-cell RNA-seq data obtained from the 10X Genomics Chromium platform was
724 processed with an in-house pipeline to identify microbial reads. This pipeline comprises

725 a series of steps designed to detect bacterial reads with high sensitivity, while discarding
726 potential false positives. For microbial identification, only the read 2 fastq file from the
727 raw sequencing files, containing the information on the cDNA fragment, was used.
728 Trimmomatic⁵⁶ (v0.38) was used to trim low quality bases and adapters, and discard
729 short reads. Additionally, Prinseq++⁵⁷ (v1.2) was used to discard reads with low-
730 complexity stretches such as poly-A sequences. Following these two quality control
731 steps, reads from human and potential sequencing artifacts (phage phiX174) were
732 mapped with STAR⁵⁸ (v2.7.1) and discarded. The remaining unmapped reads were
733 mapped against reference microbial genomes using a 2-step approach: first, we scanned
734 these remaining reads using mash screen⁵⁹ (v2.0) against a custom database of 11685
735 microbial reference genomes including bacteria, archaea, fungi and viruses. Genomes
736 likely to be present in the analyzed sample (selected using a threshold of at least two
737 shared hashes from mash screen) were selected and reads were pseudoaligned to this
738 subset of reference genomes using kallisto⁶⁰ (v0.44.0). Kallisto provides two outputs: an
739 “abundances” table containing the number of reads aligned to each gene from the pre-
740 selected set of reference genomes and a pseudo-alignment file (in *.bam format)
741 containing the mapping information for each of the reads processed by kallisto. From
742 the abundances table, we derived a taxonomy table, assigning each gene to its
743 corresponding species, as well as a functional table, mapping each gene to KEGG
744 functional annotation using KEGG Orthology numbers (KOs). To remove potential
745 artifacts, two additional filters were applied to the taxonomic table: first, if less than 10
746 different functions (i.e. 10 different KOs) were expressed from a given species, such
747 species was discarded. This filter ensures identification of active bacteria, minimizing
748 the capture of contaminants appearing during the sample preparation or sequencing.
749 Second, if one function accounted for more than 95% of the mapped reads of a given
750 species, it was also discarded. This filter was aimed at removing potential artifacts
751 caused by errors in the reference genome assemblies from our database.

752

753 Bacterial reads were assigned their specific barcodes and UMIs as follows: read IDs from
754 the mapped microbial reads were retrieved from the kallisto pseudoalignment (*.bam)
755 output using SAMtools (v1.9)⁶¹. These unique read IDs were used to retrieve the specific
756 barcodes and UMIs using the raw read 1 fastq files, thus assigning each barcode and

757 UMI univocally to a microbial species and function. Barcodes assigned to bacterial
758 species that had been removed in the last two filtering steps of the single-cell analysis
759 pipeline (see above) were discarded, to avoid including potential contaminants in the
760 host-bacteria association analyses.

761

762 Differences in lower respiratory tract microbial taxa between COVID-19 patients and
763 controls, ICU and ward patients, and invasive and non-invasive ventilation types were
764 calculated using Wilcoxon rank-sum tests on centered-log-ratio (CLR)-transformed data.
765 This more lenient approach than the one used for 16S data was chosen due to the low
766 number of samples available and the reduced number of bacterial reads identified per
767 sample. Prior to CLR data transformation, we filtered the data using the CoDaSeq.filter
768 function, to keep samples with more than 1,000 reads and taxa with a relative
769 abundance above 0.1%. Zeros were imputed using the same approach as for the 16S
770 amplicon data.

771

772 **Direct associations between bacteria and host cells**

773

774 Host single-cell transcriptomics data was obtained from the Seurat⁶² object after
775 preprocessing and integrating the samples of the single-cell cohort, as described
776 previously¹⁸. From the Seurat object, the metadata was extracted, including the
777 information on patient group (COVID-19 or control) and severity of the disease
778 (moderate or critical) as well as cell type and subtype annotation corresponding to each
779 barcode. Enrichment of bacteria detected in patient groups or cell types was calculated
780 using chi-squared tests, with effect sizes determined via the standardized residuals.
781 Significance was assessed via post-hoc tests using the R package `chisq.posthoc.test`⁶³.

782

783 To evaluate the overlap between bacterial and viral reads detection in host cells of
784 COVID-19 patients, we considered the total number of cells analyzed in these patients:
785 33,243. Of these, 31,868 cells do not have associated bacterial or viral reads; 1,032 have
786 only bacterial reads; 342 have only viral reads; and 1 has both viral and bacterial reads
787 detected (Supplementary Table 5). The marginal probability for bacterial detection is
788 thus $P(\text{bacterial detection}) = 1,033/33,243 = 0.031$; while the marginal probability for

789 viral detection in this dataset is $P(\text{viral detection}) = 343/33,243 = 0.010$. Assuming
790 independence of both events, the joint probability of finding a host cell associated to
791 both bacterial and viral reads would be $P(\text{bacterial and viral detection}) = 0.031 \cdot 0.010 =$
792 $3.2 \cdot 10^{-4}$. With this joint probability and a total of 33,243 cells profiled, an average of
793 10.65 host cells should have both bacterial and viral reads detected. A Chi-squared test
794 suggests non-independence of the data ($p\text{-value} = 4.1 \cdot 10^{-3}$). Additionally, we performed
795 an exact binomial test, considering number of successes = 1 (joint bacterial and viral
796 detection), probability of success = $3.2 \cdot 10^{-4}$, (the joint probability assuming
797 independence of both events), and number of trials = 33,243 (the total number of cells
798 studied). This two-sided test evaluates the null hypothesis that the joint probability of
799 both events is the one calculated assuming independence. The result of this test ($p\text{-value}$
800 $= 5.7 \cdot 10^{-4}$) suggests rejecting the null hypothesis. Therefore, these analyses altogether
801 suggest that both events are not independent and that there is mutual exclusion of
802 microbiome members and viruses in the same host immune cells.

803

804 For cell types showing an enrichment in associated bacteria, a new Seurat object was
805 created by subsetting the specific cell type. Chi-squared tests were also used to
806 determine enrichment of bacteria-associated cell subtypes. Previous annotations of cell
807 subtypes¹⁸ were used to generate new clusters manually and identify marker genes for
808 these subtypes, using the function `findAllMarkers` from Seurat. This function was also
809 used to find differentially expressed genes between bacteria-associated and not-
810 bacteria-associated host cells of each subtype. When using this function, reported
811 adjusted p-values are calculated using Bonferroni correction by default.

812

813 **Figure legends**

814

815 **Figure 1. Sample overview and alpha diversity.** a) Longitudinal sampling of patients.
816 Each line represents one patient. Yellow lines span the days spent in ward, while blue
817 lines span the days spent in ICU. Red points mark hospital discharge dates. Crosses
818 indicate the timepoints where swab samples were obtained for microbiome analyses.
819 b) Top 15 most abundant genera in this cohort. Samples with > 10,000 reads assigned
820 to microbial taxa at the genus level were stratified according to the sampling moment:

821 upon admission, throughout the ICU stay or at ICU discharge/during treatment in ward.
822 c) Effect of the length of ICU stay and SARS-CoV-2 viral load on upper respiratory tract
823 microbiome diversity. The plot shows the model-predicted Shannon index as a function
824 of the days in ICU, for different levels of SARS-CoV-2 viral load (selected within the range
825 of observed data). Confidence intervals for the predictions are shown in Supplementary
826 Figure 1b. d) Association of the length of ICU stay and calprotectin gene expression
827 levels with upper respiratory tract microbiome diversity. The plot shows the model-
828 predicted Shannon index as a function of the days in ICU, for different levels of
829 calprotectin (subunit S100A8) gene expression, selected within the range of observed
830 data. Confidence intervals for the predictions are shown in Supplementary Figure 1c.

831

832 **Figure 2. Upper respiratory microbiome covariates in COVID-19.** a) Significant (BH-
833 corrected p -value < 0.05) covariates explaining microbiota variation in the upper
834 respiratory tract in this cohort. Individual covariates are listed on the y-axis, their color
835 corresponds to the metadata category they belong to: technical data, disease-related,
836 microbiological tests, comorbidities or host cell populations or gene expression, the
837 latter measured with nCounter (see Methods). Darker colors refer to the individual
838 variance explained by each of these covariates assuming independency, while lighter
839 colors represent the cumulative and non-redundant variance explained by incorporating
840 each variable to a model using a stepwise dbRDA analysis (using Aitchison distances).
841 The black horizontal line separates those variables that are significant in the non-
842 redundant analysis on top (Patient ID and oxygen support type) from the rest. b) RDA
843 ordination plot showing the first 2 constrained axes. Ordination is constrained by the
844 two significant variables “Patient ID” and “Oxygen support”. Samples are depicted as
845 points, whose color indicates the oxygen support type of the patient and whose shape
846 indicates stay at ward or ICU (at the moment of sampling). Axes indicate the variance
847 explained by the first two constrained components of the RDA analysis. c) Species- (left)
848 and strain-level diversity (right) of the samples, stratified by oxygen support type. d)
849 Pearson correlation between average species- and strain- level diversity for each of the
850 oxygen support categories. e) Significant differences in taxa abundances among oxygen
851 support types. Differentially abundant taxa between invasive (red) and non-invasive
852 (blue) ventilated samples. Only the top 10 most significant taxa are shown, as

853 determined by their BH-adjusted p-value. Boxplots span from the first until the third
854 quartile of the data distribution, and the horizontal line indicates the median value of
855 the data. The whiskers extend from the quartiles until the last data point within 1.5
856 times the interquartile range, with outliers beyond. Individual data points are also
857 represented. Gray lines join samples pertaining to the same patient, taken at different
858 time points. Asterisks (*) indicate taxa that remain significant after controlling for the
859 main antibiotics (ceftriaxone and meropenem/piperacillin-tazobactam).

860

861 **Figure 3. Host single cells associated to the lower respiratory tract microbiota.** a)
862 relative proportion of cells from negative and positive COVID-19 patients with (red
863 color) and without (blue) associated bacteria. The p-value of a chi-squared test using the
864 count data is shown on top of the panel. b) Cell types enriched in bacteria-associated
865 cells. Barplots represent the proportion of cell types without (“No”) and with (“Yes”)
866 bacteria in COVID-19 positive and negative patients. For each patient class, we tested
867 for enrichment of bacteria-associated cells (“Yes”) across the different cell types, using
868 the proportions of non-bacteria associated cells (“No”) as background. Asterisks mark
869 the cell types with significant enrichment of bacteria. c) Bacterial genera preferentially
870 associating to specific cell types. The heatmaps show the standardized residuals of a chi-
871 squared test including all bacterial genera and the three host cell types enriched in
872 bacteria, for controls (left) and COVID-19 positive patients (right). Taxa with no
873 significant associations with any of the cell types are not shown. Asterisks denote
874 significant positive or negative associations: enrichments are shown in red; depletions
875 are depicted in blue. d) Host cell subtypes associated with bacteria. The heatmap shows
876 the standardized residuals of a chi-squared test including the subtypes of neutrophils,
877 monocytes and monocyte-derived macrophages with associated bacteria, considering
878 cells without bacteria as background. Asterisks denote significant positive or negative
879 associations: enrichments are shown in red; depletions are depicted in blue. e) Marker
880 genes detected for the 5 different subtypes of neutrophils. The heatmap also shows
881 within-group differences between bacteria-associated and bacteria-non-associated
882 cells. f) Myeloid cell functional gene set showing the expression of canonical pro-
883 inflammatory, anti-inflammatory and MHC genes for the two subtypes of myeloid cells
884 significantly associated with bacteria (CCL2^{hi}-macrophages and IL1B^{hi}-monocytes). The

885 heatmap also shows within group differences between bacteria-associated and
886 bacteria-non-associated cells. Statistically significant differences after multiple testing
887 correction are marked with squares. For b)-d) asterisks denote significance as follows: *
888 = p-value \leq 0.05; ** = p-value \leq 0.01; *** = p-value \leq 0.001; **** = p-value \leq 0.0001.

889

890

891 **Supplementary Figure Legends**

892

893 **Supplementary Figure 1.** Alpha diversity in the upper respiratory tract. a) Shannon
894 diversity index of all samples, stratified by the sampling moment: admission, throughout
895 ICU stay or at ICU discharge/during treatment in ward. The p-value of a Kruskal-Wallis
896 test, as well as the those of Dunn tests corresponding to the pairwise differences among
897 the three groups, are shown. b) Forest plot of the fixed effects estimates of the variables
898 selected in the best model predicting Shannon diversity index. The points and values
899 above indicate the fixed effect estimates of the variables selected, while the horizontal
900 lines span their 95% confidence intervals. Asterisks denote significance as follows: * =
901 p-value \leq 0.05; ** = p-value \leq 0.01; *** = p-value \leq 0.001; **** = p-value \leq 0.0001. c) c)
902 Effect of the length of ICU stay and SARS-CoV-2 viral load on upper respiratory tract
903 microbiome diversity. Each plot shows the model-predicted Shannon index as a function
904 of the days in ICU, for a different level of SARS-CoV-2 viral load (selected within the range
905 of observed data). Shaded areas correspond to the 95% confidence intervals. d)
906 Association of the length of ICU stay and calprotectin gene expression levels with upper
907 respiratory tract microbiome diversity. Each plot shows the model-predicted Shannon
908 index as a function of the days in ICU, for a different level of calprotectin (subunit
909 S100A8) gene expression, selected within the range of observed data. Shaded areas
910 correspond to the 95% confidence intervals. e) Model-averaged relative importance of
911 each of the variables selected (only for fixed effects). Variable importance was
912 calculated as a weighted average of the models in which each of the variables appeared,
913 with weights corresponding to the model ranks, defined by their AICc values. f) Intra-
914 patient differences of alpha-diversity values, before and after administration of
915 meropenem/piperacillin-tazobactam (left) or mechanical ventilation (right). P-values
916 shown are derived from Wilcoxon signed-rank tests. For (a,f), boxplots span from the

917 first until the third quartile of the data distribution, and the horizontal line indicates the
918 median value of the data. The whiskers extend from the quartiles until the last data
919 point within 1.5 times the interquartile range, with outliers beyond. Individual data
920 points are also represented.

921

922 **Supplementary Figure 2.** Association of antibiotics and mechanical ventilation. a)
923 Mosaic plots showing, for each category of oxygen support, the proportion of samples
924 receiving ceftriaxone administration (current administration on day of sampling, left) or
925 the proportion of samples having received meropenem or piperacillin-tazobactam
926 (ongoing or previous treatment, right). P-values denote the significance of Chi-squared
927 tests for these associations. The different oxygen support levels are: 1-oxygen flow (via
928 nasal cannula); 2-high flow oxygen support; 3-non-invasive ventilation (CPAP, BIPAP); 4-
929 invasive ventilation; 5-prone ventilation; 6-extra corporeal membrane oxygenation
930 (ECMO); 7-nitric oxide inhalation. Levels 4-7 correspond to mechanically ventilated
931 patients. b) Longitudinal sampling of patients, showing specific antibiotic
932 administration. Each line represents one patient. Yellow lines represent no-specific
933 antibiotic administration for the spanned period; blue lines represent antibiotic was
934 administered during that period. Shaded areas in light gray represent ICU stay, whilst
935 areas in dark gray represent periods with the patient receiving mechanical ventilation.
936 Crosses indicate the timepoints where swab samples were obtained for microbiome
937 analyses. Individual yellow points at later times represent follow-up visits.

938

939 **Supplementary Figure 3.** Differentially abundant taxa between oxygen support types.
940 a) The 29 taxa whose abundance is significantly different between non-invasive and
941 invasive ventilation are represented. b) The 20 taxa whose abundance is significantly
942 different between ventilation types, after controlling for antibiotic usage, are
943 represented. Boxplots span from the first until the third quartile of the data distribution,
944 and the horizontal line indicates the median value of the data. The whiskers extend from
945 the quartiles until the last data point within 1.5 times the interquartile range, with
946 outliers beyond. Individual data points are also represented. Gray lines join samples
947 pertaining to the same patient, taken at different time points.

948

949 **Supplementary Figure 4.** Absolute microbial read counts in single-cell RNA-seq data
950 from BAL samples. The top 15 species detected in our analyses are depicted. Samples
951 are grouped by disease type (control for non-COVID-19 pneumonia patients, or COVID-
952 19) and hospital stay (ICU or ward).

953

954 **Supplementary Figure 5.** Associations of specific cell types with bacteria, for COVID-19
955 and control samples. The colors represent the strength of the association as the
956 standardized residuals of a Chi-squared test. Red colors indicate a positive association
957 (i.e. enrichment) of bacteria for each cell type. Blue colors indicate a negative
958 association (i.e. depletion) of bacteria for a given cell type. Asterisks denote significance
959 as follows: * = p-value ≤ 0.05 ; ** = p-value ≤ 0.01 ; *** = p-value ≤ 0.001 ; **** = p-value
960 ≤ 0.0001 .

961

962 **Author contributions**

963 VLR, ACG, JW, JR designed the study. SJ, TVW, JN, CD, JG, GH, PM collected and
964 processed the BAL samples. PVM and LV collected the clinical data. JW and EW collected
965 the swabs. VLR, ACG and JW processed the swabs. JW, MB and SMM determined
966 SARS-CoV-2 viral loads and host gene expression from swabs. DL and JQ generated the
967 single-cell raw data as well as the processed gene-count matrix with annotations of cell
968 types and subtypes. VLR and ACG analyzed the data. VLR, ACG and JR wrote the
969 manuscript. All authors have read and approved the manuscript.

970

971 **Acknowledgments**

972 This study has been supported by funding from the VIB Grand Challenges Program. VLR
973 is supported by an FWO senior postdoctoral fellowship (12V9421N). ACG is supported
974 by an EMBO postdoctoral fellowship (ALTF 349-2019). The Raes lab is supported by KU
975 Leuven, the Rega institute and VIB.

976

977 **Conflict of interest declaration**

978 The authors declare no competing interests.

979

980 **CONTAGIOUS collaborators**

981 Yannick Van Herck, Alexander Wilmer, Michael Casaer, Stephen Rex, Nathalie Lorent,
982 Jona Yserbyt, Dries Testelmans, Karin Thevissen.

983

984 **Data availability statement**

985 Raw amplicon sequencing data that support the findings of this study have been
986 deposited at the European Genome-phenome Archive (EGA), with accession no
987 EGAS00001004951. The single cell RNA-seq data was first described in a separate
988 publication¹⁸ and deposited also in EGA with accession number EGAS00001004717.

989

990 **References**

991

- 992 1. Zhou, F. *et al.* Clinical course and risk factors for mortality of adult inpatients
993 with COVID-19 in Wuhan, China: a retrospective cohort study. *Lancet* **395**,
994 1054–1062 (2020).
- 995 2. Grasselli, G. *et al.* Risk Factors Associated with Mortality among Patients with
996 COVID-19 in Intensive Care Units in Lombardy, Italy. *JAMA Intern. Med.* **180**,
997 1345–1355 (2020).
- 998 3. Mikami, T. *et al.* Risk Factors for Mortality in Patients with COVID-19 in New
999 York City. *J. Gen. Intern. Med.* 1–10 (2020). doi:10.1007/s11606-020-05983-z
- 1000 4. Guo, W. *et al.* Diabetes is a risk factor for the progression and prognosis of
1001 <sc>COVID</sc> -19. *Diabetes. Metab. Res. Rev.* **36**, (2020).
- 1002 5. Lighter, J. *et al.* Obesity in Patients Younger Than 60 Years Is a Risk Factor for
1003 COVID-19 Hospital Admission. *Clinical infectious diseases : an official publication*
1004 *of the Infectious Diseases Society of America* **71**, 896–897 (2020).
- 1005 6. Yu, X. *et al.* SARS-CoV-2 viral load in sputum correlates with risk of COVID-19
1006 progression. *Crit. care* **24**, 170 (2020).
- 1007 7. Magleby, R. *et al.* Impact of Severe Acute Respiratory Syndrome Coronavirus 2
1008 Viral Load on Risk of Intubation and Mortality Among Hospitalized Patients With
1009 Coronavirus Disease 2019. *Clin. Infect. Dis.* (2020). doi:10.1093/cid/ciaa851
- 1010 8. Westblade, L. F. *et al.* SARS-CoV-2 Viral Load Predicts Mortality in Patients with
1011 and without Cancer Who Are Hospitalized with COVID-19. *Cancer Cell* (2020).
1012 doi:10.1016/j.ccell.2020.09.007
- 1013 9. Leisman, D. E. *et al.* Cytokine elevation in severe and critical COVID-19: a rapid
1014 systematic review, meta-analysis, and comparison with other inflammatory
1015 syndromes. *The Lancet Respiratory Medicine* **8**, 1233–1244 (2020).
- 1016 10. Sinha, P., Matthay, M. A. & Calfee, C. S. Is a ‘cytokine Storm’ Relevant to COVID-
1017 19? *JAMA Internal Medicine* **180**, 1152–1154 (2020).
- 1018 11. Khatiwada, S. & Subedi, A. Lung microbiome and coronavirus disease 2019
1019 (COVID-19): Possible link and implications. *Human Microbiome Journal* **17**,
1020 100073 (2020).
- 1021 12. Dickson, R. P., Martinez, F. J. & Huffnagle, G. B. The role of the microbiome in
1022 exacerbations of chronic lung diseases. *The Lancet* **384**, 691–702 (2014).
- 1023 13. Huffnagle, G. B., Dickson, R. P. & Lukacs, N. W. The respiratory tract microbiome

- 1024 and lung inflammation: A two-way street. *Mucosal Immunology* **10**, 299–306
1025 (2017).
- 1026 14. Zhang, H. *et al.* Metatranscriptomic Characterization of COVID-19 Identified A
1027 Host Transcriptional Classifier Associated With Immune Signaling. *Clin. Infect.*
1028 *Dis.* (2020). doi:10.1093/cid/ciaa663
- 1029 15. Xu, R. *et al.* Temporal dynamics of human respiratory and gut microbiomes
1030 during the course of COVID. *medRxiv* 2020.07.21.20158758 (2020).
1031 doi:10.1101/2020.07.21.20158758
- 1032 16. Han, Y., Jia, Z., Shi, J., Wang, W. & He, K. The active lung microbiota landscape of
1033 COVID-19 patients. *medRxiv* 2020.08.20.20144014 (2020).
1034 doi:10.1101/2020.08.20.20144014
- 1035 17. Mostafa, H. H. *et al.* Metagenomic Next-Generation Sequencing of
1036 Nasopharyngeal Specimens Collected from Confirmed and Suspect COVID-19
1037 Patients Downloaded from. *MBio* **11**, (2020).
- 1038 18. Wauters, E. *et al.* Discriminating mild from critical COVID-19 by innate and
1039 adaptive immune single-cell profiling of bronchoalveolar lavages. *Cell Res.* **31**,
1040 272–290 (2021).
- 1041 19. Ho Man, W., de Steenhuijsen Pitters, W. A. & Bogaert, D. The microbiota of the
1042 respiratory tract: gatekeeper to respiratory health. (2017).
1043 doi:10.1038/nrmicro.2017.14
- 1044 20. Whelan, F. J. *et al.* Longitudinal sampling of the lung microbiota in individuals
1045 with cystic fibrosis. *PLoS One* **12**, (2017).
- 1046 21. Nolan, T. J. *et al.* Low-pathogenicity *Mycoplasma* spp. alter human monocyte
1047 and macrophage function and are highly prevalent among patients with
1048 ventilator-acquired pneumonia. *Thorax* **71**, 594–600 (2016).
- 1049 22. Harris, J. R., Balajee, S. A. & Park, B. J. Pneumocystis jirovecii pneumonia:
1050 Current knowledge and outstanding public health issues. *Current Fungal*
1051 *Infection Reports* **4**, 229–237 (2010).
- 1052 23. Dickson, R. P. *et al.* Cell-associated bacteria in the human lung microbiome.
1053 *Microbiome* **2**, 28 (2014).
- 1054 24. Wang, S. *et al.* S100A8/A9 in inflammation. *Frontiers in Immunology* **9**, 1298
1055 (2018).
- 1056 25. Coveney, A. P. *et al.* Myeloid-related protein 8 induces self-tolerance and cross-
1057 tolerance to bacterial infection via TLR4- and TLR2-mediated signal pathways.
1058 *Nat. Publ. Gr.* (2015). doi:10.1038/srep13694
- 1059 26. Wauters, E. *et al.* Discriminating Mild from Critical COVID-19 by Innate and
1060 Adaptive Immune Single-cell Profiling of Bronchoalveolar Lavages. *Patrick*
1061 *Matthys* **9**, 14 (2020).
- 1062 27. Holm, J. & Hansen, S. I. Characterization of soluble folate receptors (folate
1063 binding proteins) in humans. Biological roles and clinical potentials in infection
1064 and malignancy. *Biochimica et Biophysica Acta - Proteins and Proteomics* **1868**,
1065 140466 (2020).
- 1066 28. Zakharkina, T. *et al.* The dynamics of the pulmonary microbiome during
1067 mechanical ventilation in the intensive care unit and the association with
1068 occurrence of pneumonia. *Thorax* **72**, 803–810 (2017).
- 1069 29. Schmitt, F. C. F. *et al.* Pulmonary microbiome patterns correlate with the course
1070 of disease in patients with sepsis-induced ARDS following major abdominal

- 1071 surgery. (2020). doi:10.1016/j.jhin.2020.04.028
- 1072 30. Kelly, B. J. *et al.* Composition and dynamics of the respiratory tract microbiome
1073 in intubated patients. *Microbiome* **4**, 7 (2016).
- 1074 31. Otsuji, K. *et al.* Dynamics of microbiota during mechanical ventilation in
1075 aspiration pneumonia. *BMC Pulm. Med.* **19**, 260 (2019).
- 1076 32. Shen, Z. *et al.* Genomic Diversity of Severe Acute Respiratory Syndrome-
1077 Coronavirus 2 in Patients With Coronavirus Disease 2019. *Clinical infectious*
1078 *diseases : an official publication of the Infectious Diseases Society of America* **71**,
1079 713–720 (2020).
- 1080 33. Marotz, C. *et al.* Title: Microbial context predicts SARS-CoV-2 prevalence in
1081 patients and the hospital built environment. *medRxiv* 2020.11.19.20234229
1082 (2020). doi:10.1101/2020.11.19.20234229
- 1083 34. Bassis, C. M. *et al.* Analysis of the upper respiratory tract microbiotas as the
1084 source of the lung and gastric microbiotas in healthy individuals. *MBio* **6**, (2015).
- 1085 35. Brennan, M. T. *et al.* The role of oral microbial colonization in ventilator-
1086 associated pneumonia. *Oral Surgery, Oral Med. Oral Pathol. Oral Radiol.*
1087 *Endodontology* **98**, 665–672 (2004).
- 1088 36. Stonecypher, K. Ventilator-Associated Pneumonia: The Importance of Oral Care
1089 in Intubated Adults. *Crit. Care Nurs. Q.* **33**, 339–347 (2010).
- 1090 37. Segal, L. N. *et al.* Enrichment of the lung microbiome with oral taxa is associated
1091 with lung inflammation of a Th17 phenotype. *Nat. Microbiol.* **1**, 1–11 (2016).
- 1092 38. Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K. & Schloss, P. D.
1093 Development of a dual-index sequencing strategy and curation pipeline for
1094 analyzing amplicon sequence data on the miseq illumina sequencing platform.
1095 *Appl. Environ. Microbiol.* **79**, 5112–5120 (2013).
- 1096 39. Moens, B. *et al.* Simultaneous RNA quantification of human and retroviral
1097 genomes reveals intact interferon signaling in HTLV-1-infected CD4+ T cell lines.
1098 *Virology* **9**, 171 (2012).
- 1099 40. Fukutani, K. F. *et al.* Pathogen transcriptional profile in nasopharyngeal aspirates
1100 of children with acute respiratory tract infection. *J. Clin. Virol.* **69**, 190–196
1101 (2015).
- 1102 41. Bouzas, M. L. *et al.* Diagnostic accuracy of digital RNA quantification versus real-
1103 time PCR for the detection of respiratory syncytial virus in nasopharyngeal
1104 aspirates from children with acute respiratory infection. *J. Clin. Virol.* **106**, 34–40
1105 (2018).
- 1106 42. Fukutani, K. F. *et al.* In situ immune signatures and microbial load at the
1107 nasopharyngeal interface in children with acute respiratory infection. *Front.*
1108 *Microbiol.* **9**, (2018).
- 1109 43. Hildebrand, F., Tadeo, R., Voigt, A. Y., Bork, P. & Raes, J. LotuS: An efficient and
1110 user-friendly OTU processing pipeline. *Microbiome* **2**, 30 (2014).
- 1111 44. Callahan, B. J. *et al.* DADA2: High-resolution sample inference from Illumina
1112 amplicon data. *Nat. Methods* **13**, 581–583 (2016).
- 1113 45. McMurdie, P. J. & Holmes, S. phyloseq: An R Package for Reproducible
1114 Interactive Analysis and Graphics of Microbiome Census Data. *PLoS One* **8**,
1115 e61217 (2013).
- 1116 46. Davis, N. M., Proctor, D. M., Holmes, S. P., Relman, D. A. & Callahan, B. J. Simple
1117 statistical identification and removal of contaminant sequences in marker-gene

- 1118 and metagenomics data. *Microbiome* **6**, 226 (2018).
- 1119 47. Oksanen, J. *et al.* vegan: Community Ecology Package. (2019).
- 1120 48. Gloor, G. B., Wu, J. R., Pawlowsky-Glahn, V. & Egozcue, J. J. It's all relative:
1121 analyzing microbiome data as compositions. *Ann. Epidemiol.* **26**, 322–329
1122 (2016).
- 1123 49. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and
1124 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 1125 50. Pagès, H., Aboyoun, P., Gentleman, R. & DebRoy, S. Biostrings: Efficient
1126 manipulation of biological strings. (2019).
- 1127 51. Kassambara, A. rstatix: Pipe-Friendly Framework for Basic Statistical Tests.
1128 (2020).
- 1129 52. Calcagno, V. & de Mazancourt, C. glmulti: An R package for easy automated
1130 model selection with (generalized) linear models. *J. Stat. Softw.* **34**, 29 (2010).
- 1131 53. Lüdtke, D. sjPlot: Data Visualization for Statistics in Social Science. (2021).
- 1132 54. Wright, E. S. *Using DECIPHER v2.0 to Analyze Big Biological Sequence Data in R.*
- 1133 55. Gregory, A. C. *et al.* Marine DNA Viral Macro- and Microdiversity from Pole to
1134 Pole. *Cell* **177**, 1109–1123.e14 (2019).
- 1135 56. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for
1136 Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 1137 57. Cantu, V. A., Sadural, J. & Edwards, R. PRINSEQ++, a multi-threaded tool for fast
1138 1 and efficient quality control and 2 preprocessing of sequencing datasets.
1139 (2019). doi:10.7287/peerj.preprints.27553v1
- 1140 58. Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–
1141 21 (2013).
- 1142 59. Ondov, B. D. *et al.* Mash Screen: high-throughput sequence containment
1143 estimation for genome discovery. *Genome Biol.* **20**, 232 (2019).
- 1144 60. Bray, N. L., Pimentel, H., Melsted, P. & Pachter, L. Near-optimal probabilistic
1145 RNA-seq quantification. *Nat. Biotechnol.* **34**, 525–527 (2016).
- 1146 61. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
1147 **25**, 2078–2079 (2009).
- 1148 62. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-cell
1149 transcriptomic data across different conditions, technologies, and species. *Nat.*
1150 *Biotechnol.* **36**, 411–420 (2018).
- 1151 63. Ebbert, D. chisq.posthoc.test: A Post Hoc Analysis for Pearson's Chi-Squared
1152 Test for Count Data. (2019).
- 1153

Figure 1

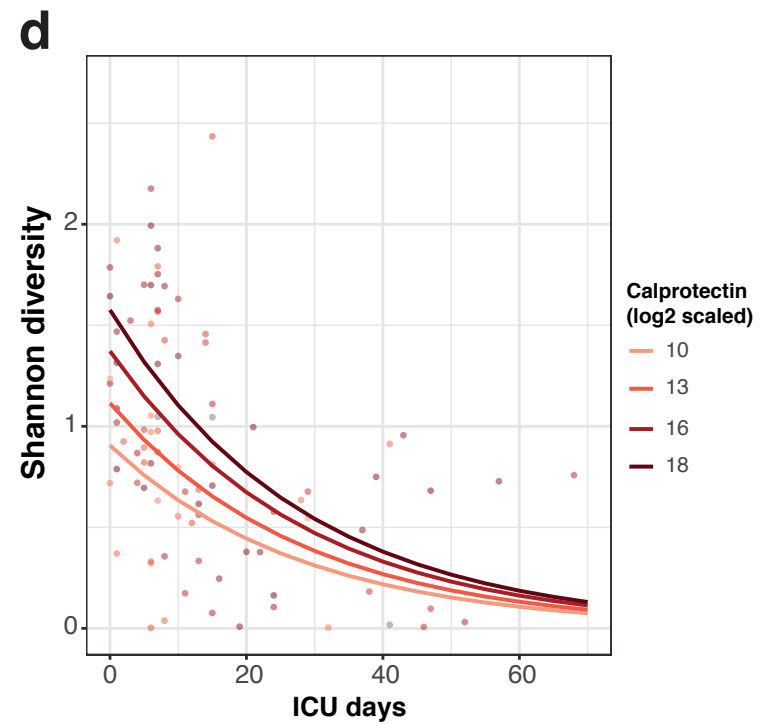
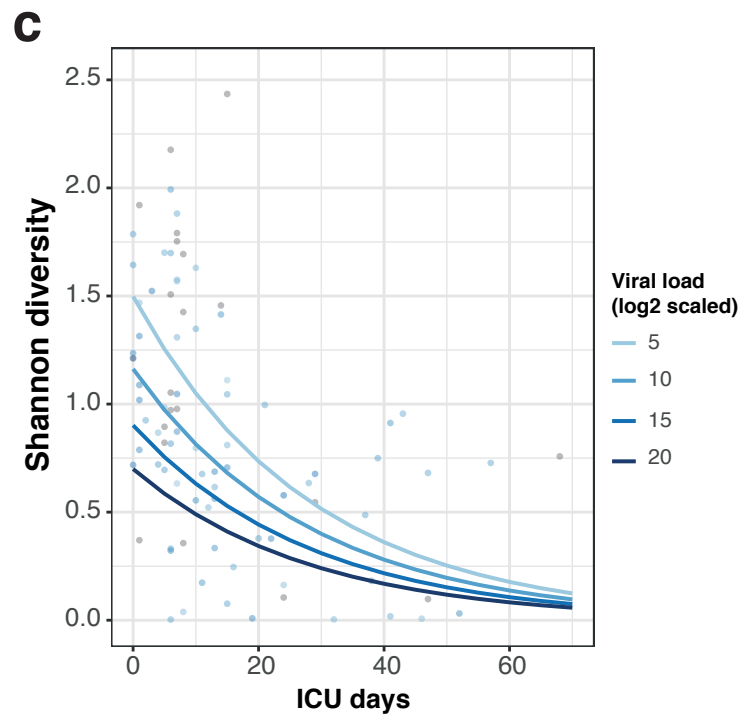
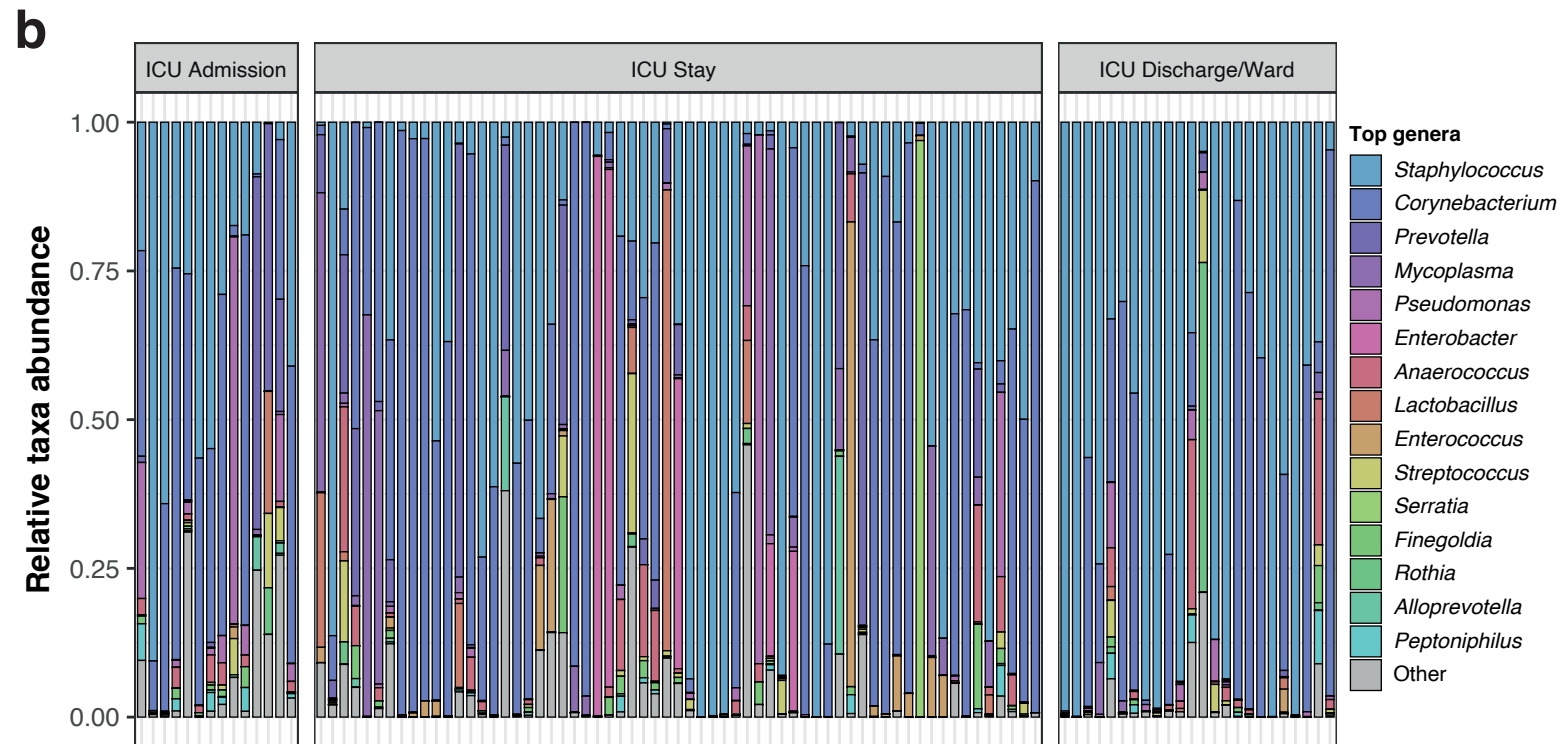
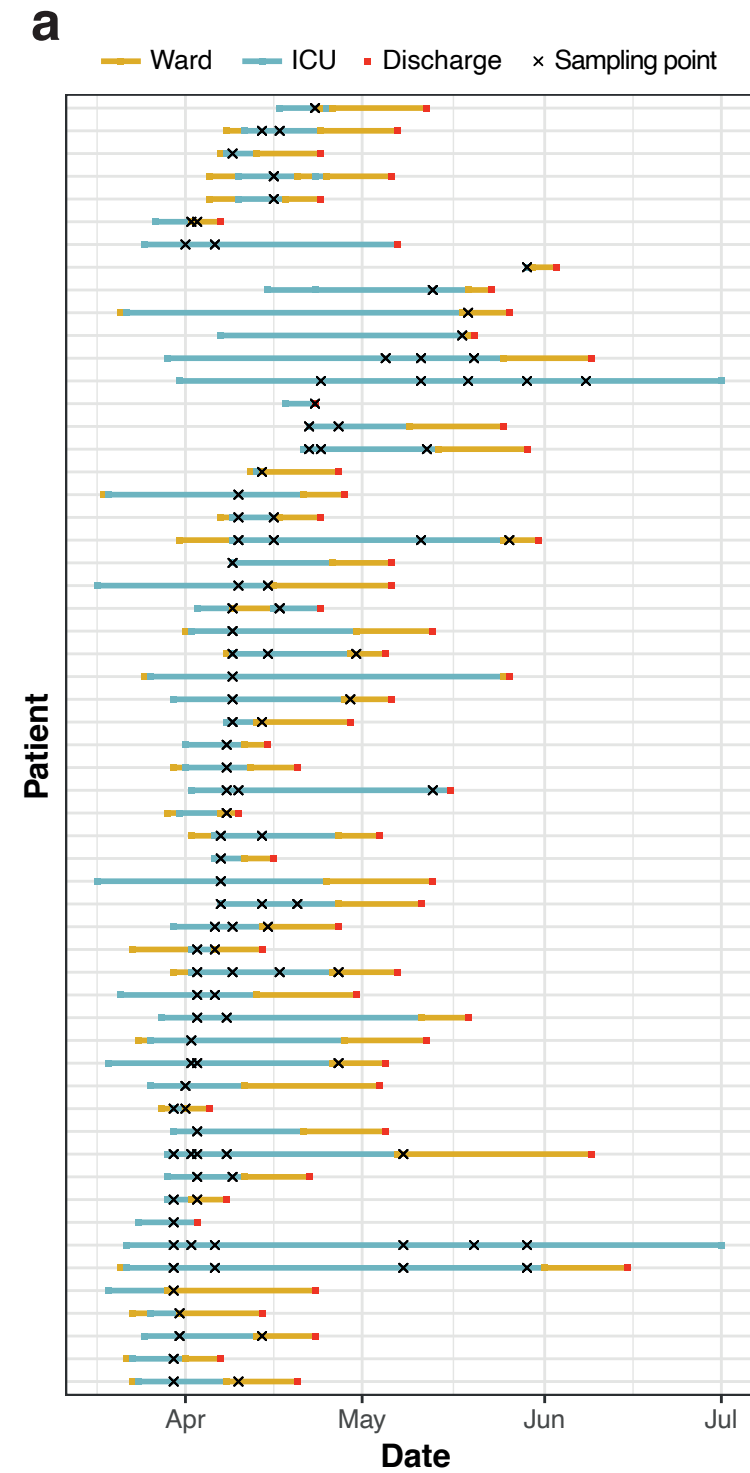
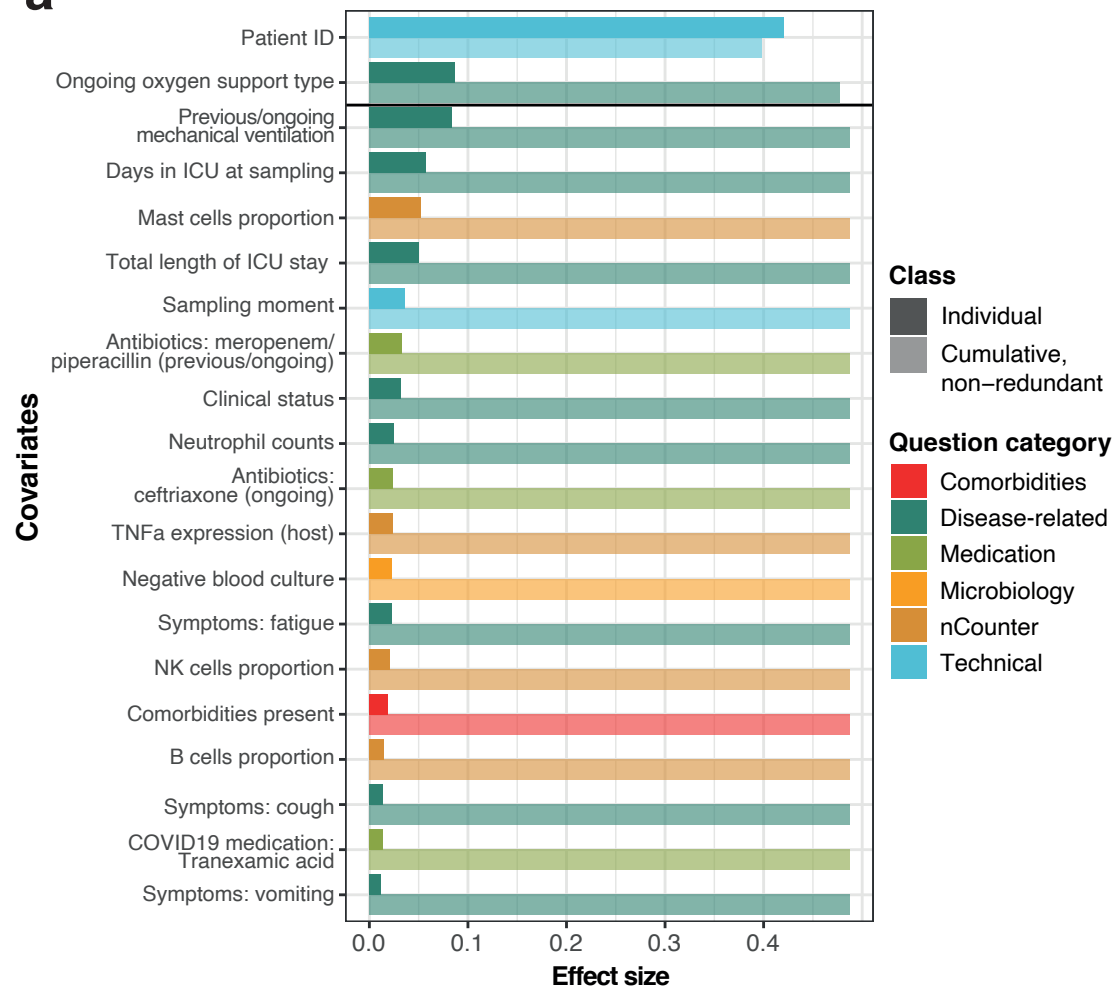
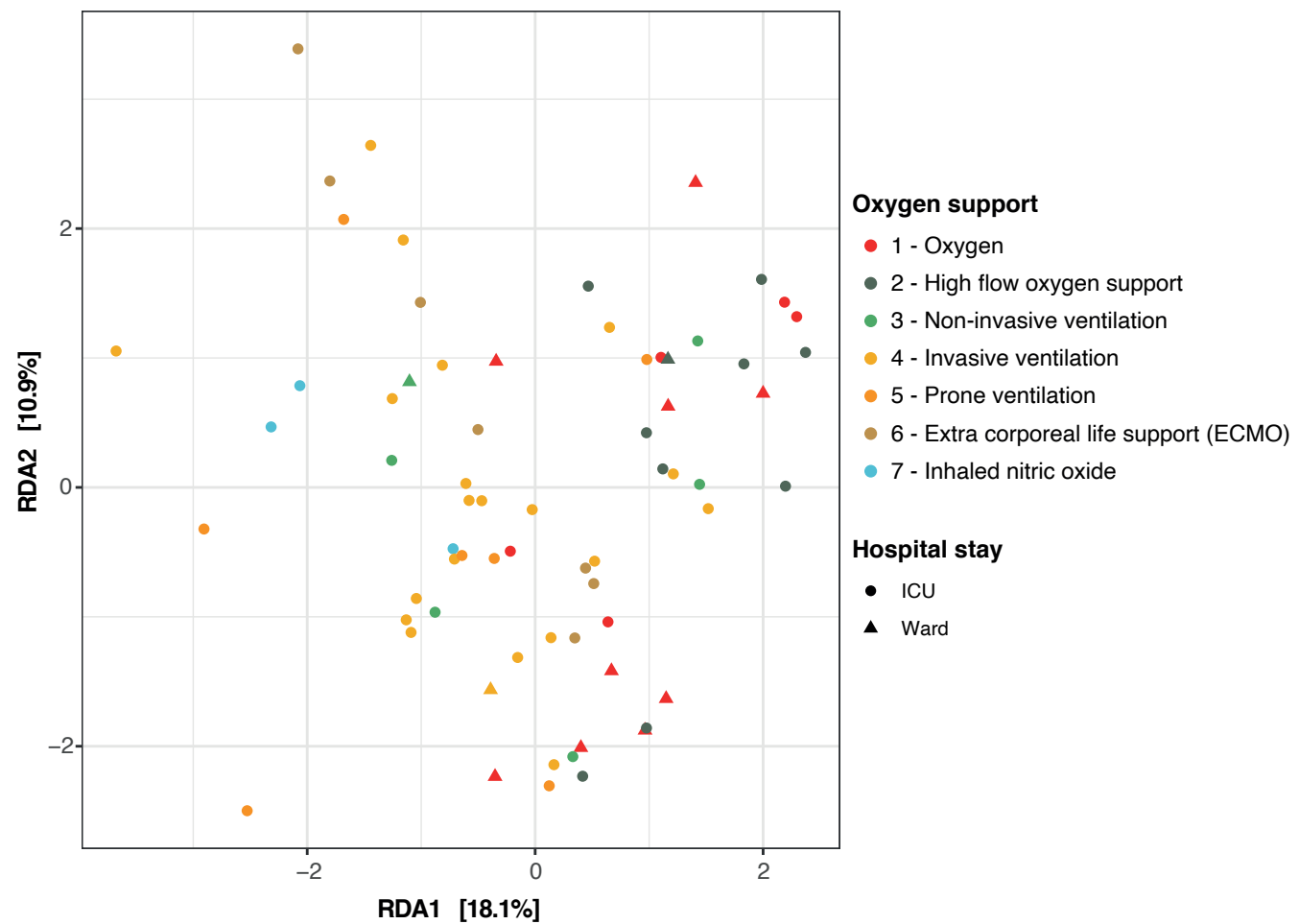


Figure 2

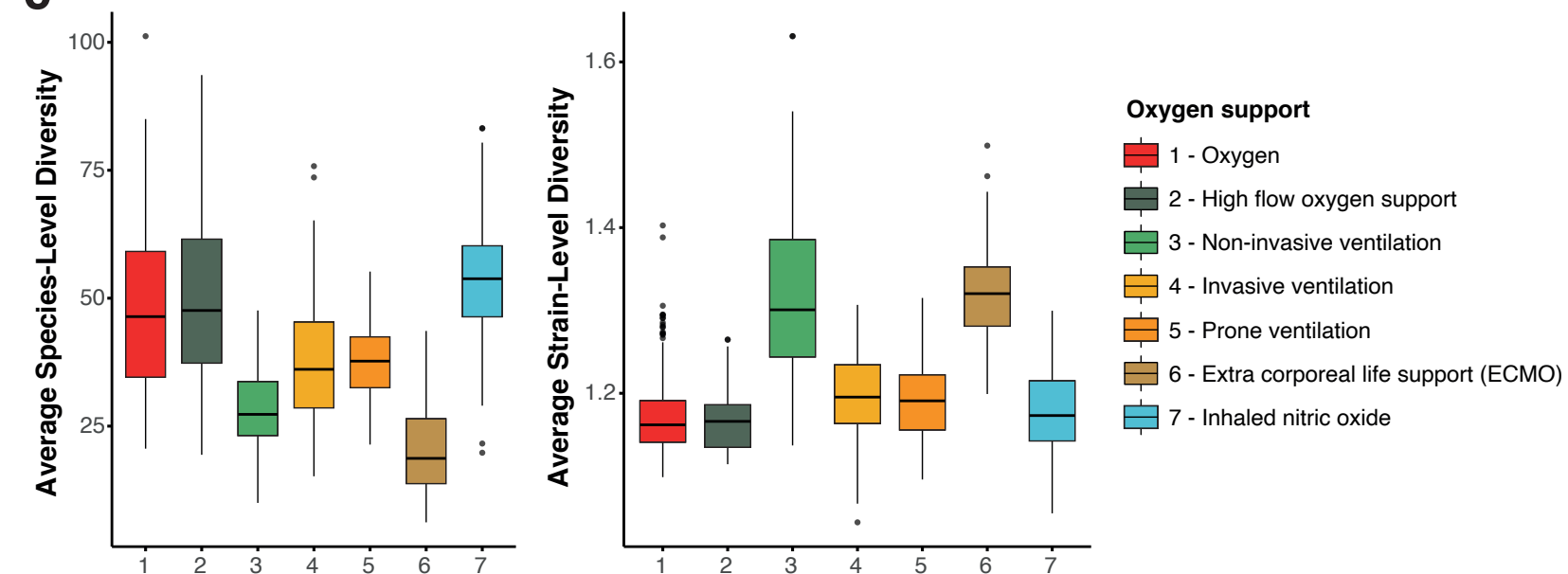
a



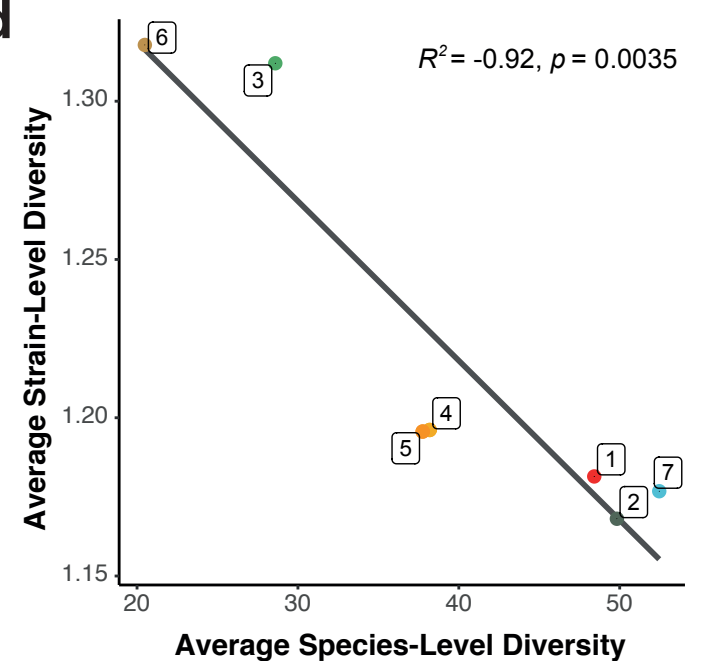
b



c



d



e

Significant differences among oxygen support types (top 10 taxa)

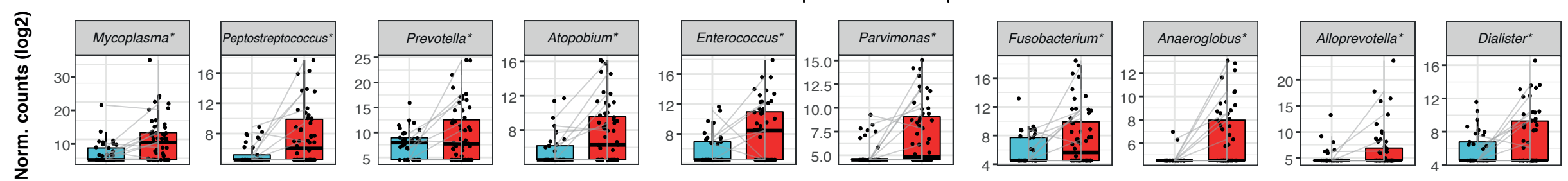
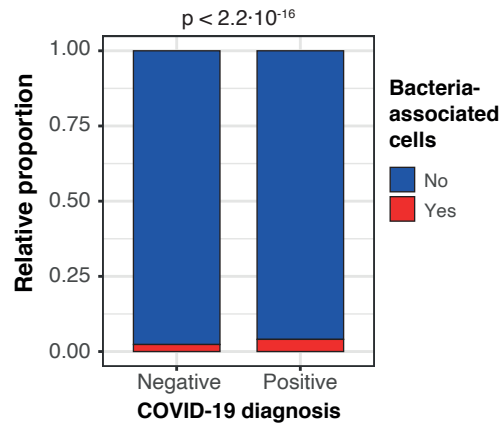
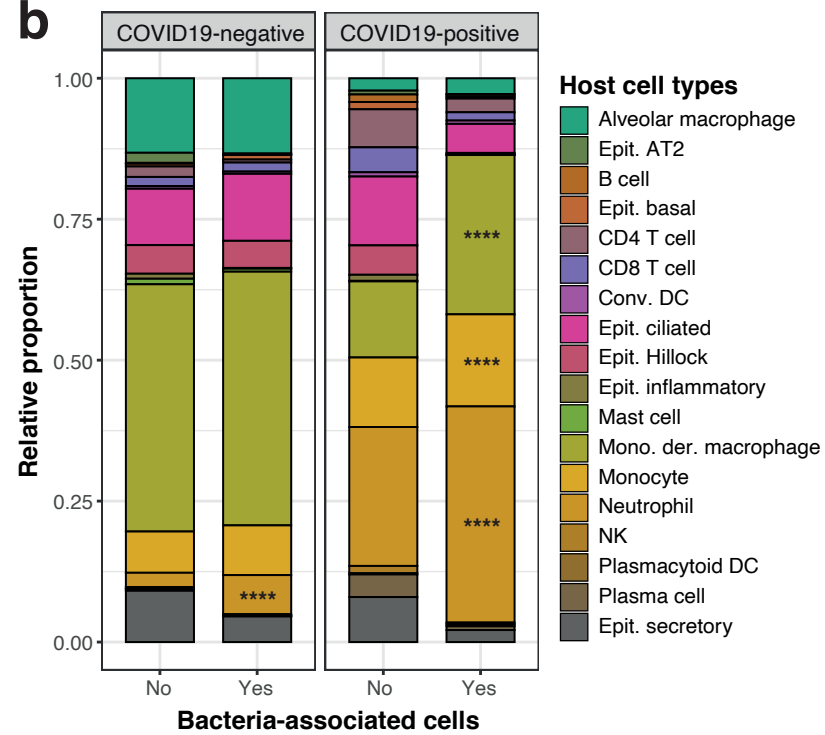


Figure 3

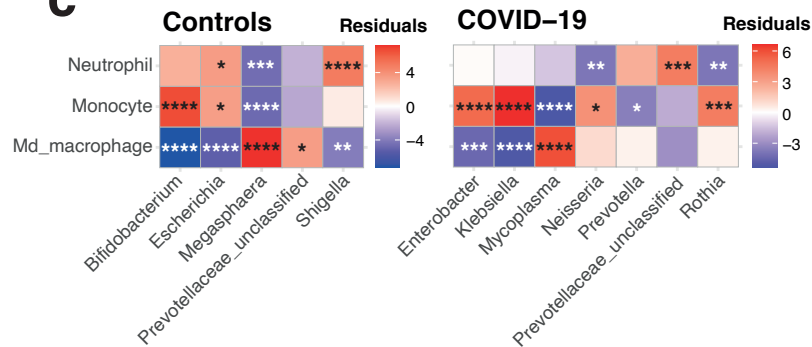
a



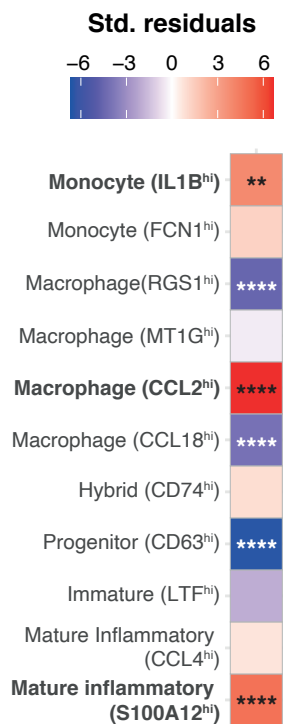
b



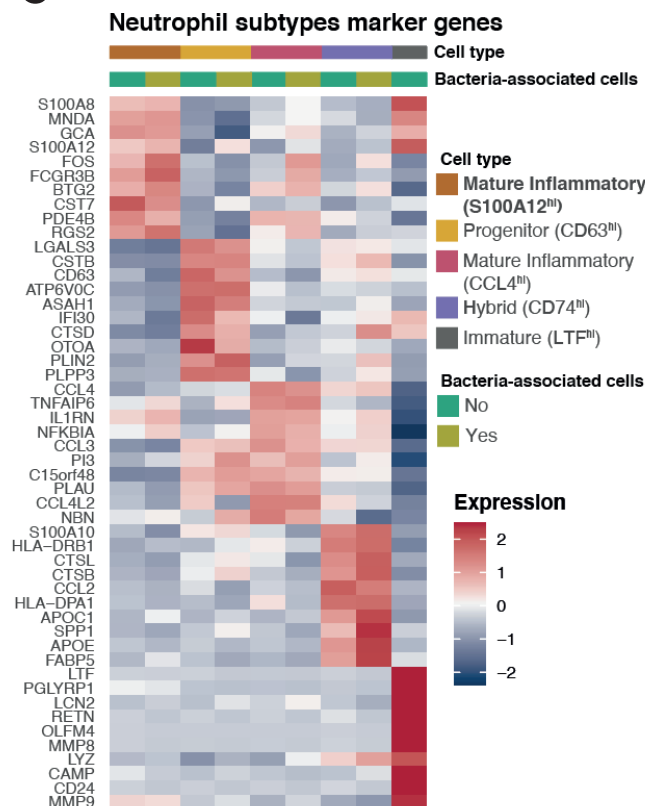
c



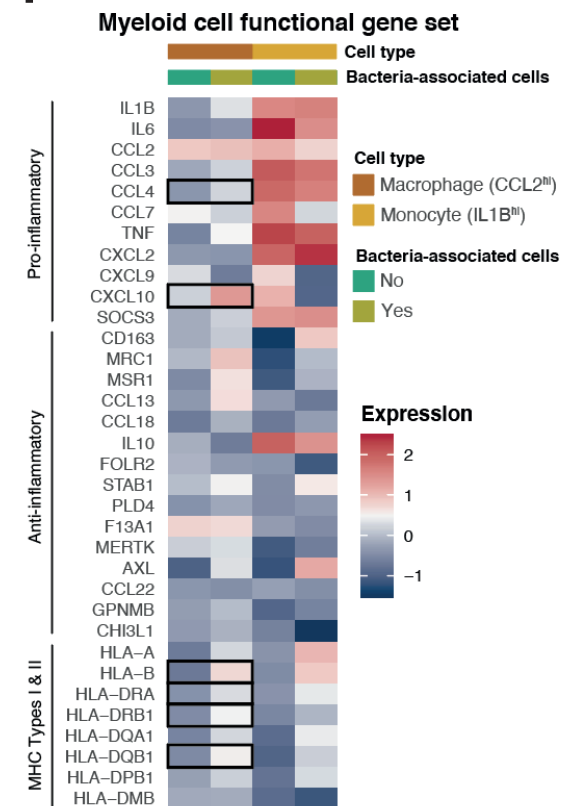
d



e



f



	<i>Upper respiratory tract (swabs)</i>	<i>Lower respiratory tract (BAL)</i>
<i>Number of patients</i>	58	35
<i>Type of sampling</i>	Longitudinal	Cross-sectional
<i>COVID-19 diagnosis (%)</i>	58 (100%)	22 (63%)
<i>Patients admitted to ICU (%)</i>	58 (100%)	21 (60%) at sampling
<i>Age (range)</i>	61 (37-83)	64 (45-85)
<i>Female sex (%)</i>	13 (22%)	12 (34%)
<i>BMI (range)</i>	29 (22-47)	26 (16-36)
<i>Diabetic (%)</i>	12 (21%)	6 (17%)
<i>Days in ICU (range)</i>	21.4 (2-72)	Not Available (cross-sectional cohort)
<i>Days in hospital (range)</i>	32.5 (6-86)	Not Available (cross-sectional cohort)

Table 1. Patient demographics of upper and lower respiratory tract cohorts