# Nanopore 16S rRNA sequencing reveals alterations in nasopharyngeal microbiome and enrichment of *Mycobacterium* and *Mycoplasma* in patients with COVID 19

Soumendu Mahapatra<sup>1\*</sup>, Rasmita Mishra<sup>1\*</sup>, Punit Prasad<sup>1\*\$</sup>, Krushna Chandra Murmu<sup>1</sup>, Shifu

Aggarwal<sup>1</sup>, Manisha Sethi<sup>1</sup>, Priyanka Mohapatra<sup>1</sup>, Arup Ghosh<sup>1</sup>, Rina Yadav<sup>1</sup>, Hiren Dodia<sup>1</sup>,

Shamima Azma Ansari<sup>1</sup>, Saikat De<sup>1</sup>, Deepak Singh<sup>1</sup>, Amol Suryawanshi<sup>1</sup>, Rupesh Dash<sup>1</sup>,

Shantibhushan Senapati<sup>1</sup>, Tushar K. Beuria<sup>1</sup>, Soma Chattopadhyay<sup>1</sup>, Gulam

Hussain Syed<sup>1</sup>, Rajeeb Swain<sup>1</sup>, Sunil K. Raghav<sup>1</sup>, Ajay Parida<sup>1\$</sup>

<sup>1</sup>Institute of Life Sciences, Bhubaneswar, Odisha, India.

\* These authors contributed equally to this work

<sup>\$</sup> Correspondence: Ajay Parida, Ph.D.

Institute of Life Sciences, Nalco Square, Chandrasekharpur, Bhubaneswar, Odisha – 751023 Phone: +91-674-2304324 Email: ajayparida@ils.res.in drajayparida@gmail.com

#### Punit Prasad, Ph.D.

Institute of Life Sciences, Nalco Square, Chandrasekharpur, Bhubaneswar, Odisha – 751023 Phone: +91-674-2304319 Email: <u>punit@ils.res.in punit.ils@gov.in</u>

**Running title:** 16S rRNA sequencing of nasopharyngeal microbiome using Oxford Nanopore<sup>TM</sup>

Keywords: Nasopharyngeal microbiome, Nanopore, COVID-19, 16S rRNA

#### 1 Abstract

2 The coronavirus disease 2019 (COVID-19) pandemic caused by severe acute respiratory syndrome corona virus 2 (SARS-CoV-2) is a major global health concern. This virus infects 3 4 the upper respiratory tract and causes pneumonia-like symptoms. So far, few studies have shown that respiratory infections alter nasopharyngeal (NP) microbiome diversity and enrich 5 opportunistic pathogens. In this study, we have sequenced the 16S rRNA variable regions, V1 6 7 through V9, extracted from NP samples of control and COVID-19 (symptomatic and 8 asymptomatic) participants using the Oxford Nanopore<sup>TM</sup> technology. Comprehensive 9 bioinformatics analysis investigating the alpha/beta diversities, non-metric multidimensional scaling, correlation studies, canonical correspondence analysis, linear discriminate analysis, 10 and dysbiosis index analysis revealed control and COVID-19-specific NP microbiomes. We 11 12 observed significant dysbiosis in COVID-19 NP microbiome with abundance of opportunistic pathogens such as Cutibacterium, Corynebacterium, Oerskovia, and Cellulomonas in 13 asymptomatic patients, and of *Streptomyces* and *Mycobacteriaceae* family in symptomatic 14 15 patients. Furthermore, we observed sharp rise in enrichment of opportunistic pathogens in symptomatic patients, with abundance of *Mycobacteria* and *Mycoplasma*, which strongly 16 17 correlated with the occurrences of chest pain and fever. Our findings contribute novel insights regarding emergence of opportunistic pathogens in COVID-19 patients and their relationship 18 with symptoms, suggesting their potential role in coinfections. 19

- 20
- 21
- 22
- 23

#### 24 Introduction

The coronavirus disease 2019 (COVID-19) pandemic, a global health threat, is caused by 25 severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2). The symptoms range from 26 27 fever, throat pain, loss of taste and smell to severe congestion in the chest, drop in oxygen levels, pneumonia, and acute respiratory distress syndrome (1). Furthermore, a significant 28 population worldwide remains asymptomatic, which is considered spreaders of the infection 29 (2). The virus enters the host via the upper respiratory tract (URT) where the spike protein 30 binds to the angiotensin I converting enzyme 2 (ACE2) receptor, an essential step in invading 31 host cells to cause progressive disease (3, 4). Random mutations in the SARS-CoV2 spike 32 protein and receptor-binding domain promote efficient invasion and enhance pathogenicity 33 (5). 34

35 The nasopharyngeal tract is inhabited by a large number of microbial communities which maintain normal homeostasis (6). Studies have revealed association between microbial 36 communities that influence viral infections of the lung, such as chronic rhinosinusitis, asthma, 37 pneumonia, and cystic fibrosis in the URT (7, 8). URT microbiome dysbiosis may also 38 enhance the opportunistic pathogen population and promote coinfection in the host (9, 10). 39 Reports have shown that nasopharyngeal (NP) swabs in viral transport media can be used to 40 investigate the NP microbial composition in patients with COVID-19 (11, 12). Recent studies 41 have revealed overall compositional changes in the NP microbiota and promotion of 42 43 opportunistic pathogens such as *Rothia* and *Veillonella* in COVID-19 patients with shortness of breath (11, 13, 14). The secondary infection in patients with COVID-19 is associated with 44 abundance of opportunistic pathogens such as Moraxella, Corynebacterium, Haemophilus, 45 Stenotrophomonas, Acinetobacter, Fusobacterium periodonticum, and Pseudomonas 46

47 aeruginosa (15-18). Studies on functional pathways of the NP metagenomics have revealed that the abundance of NP commensal bacteria such as Gemella morbillorum, Gemella 48 haemolysans, and Leptotrichia hofstadii was reduced in the respiratory tract of COVID-19 49 50 patients, indicating the role of distinct functional metabolic pathways in this infection (19, 20). Little is known about the crosstalk between SARS-CoV-2 viral infection and NP microbiota. 51 52 Moreover, systematic data connecting COVID-19-associated symptoms with microbial composition is lacking. The absence of an animal model makes it difficult to test and validate 53 the role of NP microbiota in SARS-CoV-2 infection. Studies so far have shown differences in 54 55 the abundance of different opportunistic pathogens in the NP microbiota of patients, which is one of the bottlenecks in this area of research. Hence more studies on the NP microbiome are 56 required for understanding its role in symptomatic and asymptomatic COVID-19 patients and 57 its relation with symptom severity. 58

59 In this study, we have investigated the alterations in the NP microbial ecosystem of patients with active COVID-19 (n = 46) and compared them with that of healthy individuals (n = 12). 60 We have used the 16S metagenome approach and long-read sequencing (V1–V9) with the 61 62 Nanopore sequencing method to elucidate the reduction in microbial diversity in patients with COVID-19. The composition of the NP microbiota changed significantly between 63 symptomatic and asymptomatic patients, resulting in enrichment of opportunistic pathogens 64 Interestingly, we found abundance of *Mycoplasma* and *Mycobacterium* at the genus level, 65 which strongly correlated with chest pain and fever in the symptomatic patients. 66

#### 67 Materials and Methods

Ethical approval: Ethical permission for nasopharyngeal microbiome study and the
biorepository was obtained from the Institutional Ethical Committee (IEC)/Institutional

Review Board (IRB) of the Institute of Life Sciences [(102/HEC/2020) and (100/HEC/2020)].
Approval was also obtained from the Institutional Biosafety Committee (IBSC) (V-122MISC/2007-08/01/2/2.1) for this study and the biorepository (V-122-MISC/2007-08/01) and
from the Review Committee on Genetic Manipulations (RCGM) under Department of
Biotechnology, Ministry of Science and Technology.

Sample collection and reverse transcription-polymerase chain reaction (RT-PCR): In 75 total, 60 NP samples were collected for 16S rDNA amplicon sequencing from the Institute of 76 Life Science (ILS) COVID-19 sample biorepository unit. The COVID-19-positive samples (n 77 78 = 47) were confirmed by amplifying the genes encoding SARS-CoV-2 nucleocapsid, spike, 79 and ORF1ab/RdRP using either TaqPath<sup>™</sup> COVID-19 combo kit (Invitrogen, A47814) or Meril COVID-19 one-step RT-PCR kit (Meril Diagnostics, NCVPCR-02). All samples were 80 collected in the hospital setup prior to the medication. These COVID-19-positive patients were 81 not treated with antibiotics as the patients were not aware of their COVID-19 testing results. 82 The COVID-19 patients were grouped as symptomatic (n = 22) or asymptomatic (n = 25)83 based on their clinical data. The control samples (n = 13) were negative for SARS-CoV-2 84 virus RNA and none of the subjects from whom the samples were obtained had any flu-like 85 86 symptoms. All samples were collected in viral transport media (VTM) and stored at -80°C until DNA isolation. 87

DNA extraction and PCR amplification: DNA was isolated using the PureLink<sup>™</sup>
microbiome DNA purification kit (Invitrogen, A29790) according to the manufacturer's
protocol and eluted in 40 µl elution buffer. The quality and quantity of DNA were determined
using the Multiskan<sup>™</sup>GO spectrophotometer (Thermo Scientific). 16S rDNA amplification,
library preparation, and sequencing: V1-V9 variable regions of the 16S rRNA gene were

93	amplified using 130-F (5'-GGCGGATCCAAGGAGGTGTTCCAGCCGC-3') and 139-R (5'-
94	GGCCTCGAGAGAGTTTGATCCTGGCTCAGG-3') primers. PCR (50 µl) was set up using
95	total DNA (10 ng) isolated from NP samples, primers (5 nM), and NEB Q5® High-Fidelity 2
96	X master mix (NEB, M0492L) per the manufacturer's protocol. The amplicons (~1.6 kb) were
97	analyzed on 0.8% agarose gel and cleaned using DNA Clean and Concentrator-25 kit (Zymo
98	Research, D4034). The PCR products were quantified using a Qubit 4 fluorometer (Thermo
99	Scientific) using the Qubit® dsDNA BR assay kit (Thermo Scientific, Q32853). Amplicon
100	libraries were generated following the Oxford Nanopore 1D library preparation protocol using
101	the PCR barcoding (96) genomic DNA kit (Oxford Nanopore <sup>™</sup> , SQK-LSK109). Equimolar
102	amounts of amplicon libraries were pooled and sequenced using the MinION OXFORD
103	NANOPORE <sup>TM</sup> device at the ILS DNA sequencing facility.
104	Microbiome data processing: RAW fast5 files were generated using the MinKNOW <sup>™</sup> tool
105	for individual samples. Base calling was performed using the Guppy base-caller and fastq files
106	were generated. FastQC of each sample was performed using the Babraham fastqc suite
107	(https://www.bioinformatics.babraham.ac.uk/projects/fastqc/), followed by trimming of low
108	quality reads using nanoflit. Operational taxonomic units (OTUs) were generated using
109	Kraken2 (https://ccb.jhu.edu/software/kraken2/index.shtml) (21) and the unclassified reads
110	were filtered for downstream analysis using the 'phyloseq' 'R" package to generate combined
111	OTUs for all the samples and metadata (Supplemental Table 1). Read counts for mitochondria
112	and chloroplast were discarded. Normalization and differential OTU abundance were
113	determined between control, and symptomatic and asymptomatic subjects using the DESeq2
114	function (cutoff of p-value $\leq$ 0.05). The accession ID in NCBI is PRJNA774098.

115 In-depth microbiome data analysis:

116 Diversity analysis: Alpha diversity was assessed using the Shannon diversity index and 117 Simpson Diversity index. Statistical significance was estimated using the Wilcoxon rank sum test. The beta diversity significance among groups was examined with PERMANOVA (p-118 119 value 0.001). Ordination analysis was performed by PCoA, NMDS and CCA. R packages used are 'microbiome', 'Vegan', 'ade4', ggpubr for analysis and 'ggplot2' for visualization. 120 Dysbiosis index: Microbiome dysbiosis in each sample was calculated based on Bray-Curtis 121 distances. All samples were subjected to PCoA using Bray-Curtis distances. Next, the centroid 122 (median) of the control subjects was calculated along PCoA axes. The dysbiosis score for each 123 124 sample was calculated as a Euclidian distance between its position in the PCoA space and control centroid (DI(X, HC) =  $\sqrt{(Xi - HCi)^2 + (Xj - HCj)^2}$  | (DI: Dysbiosis Index, X: Samples, HC: 125 Control Centroid). Their significance was assessed using Wilcoxon and Kruskal-Wallis test 126 (22). 127

<u>Sample correlation:</u> Correlation matrix between samples and OTUs for each taxonomic level
 (phylum, order, family, and genus) from differential OTUs was obtained using Spearman's
 correlation method and it was visualized as a heat map. Correlation coefficients for each
 sample correlation pair and each classification level and density plot were plotted with mean
 and median. The Kolmogorov test (KS) was used to determine the significance in sample
 groups (control, asymptomatic, and symptomatic).

Linear discriminant analysis (LDA) effect size (LEfSe) analysis: The LEfSe was calculated using the online Galaxy web application with the Huttenhower lab's tool (https://huttenhower.sph.harvard.edu/galaxy/). LDA effect size was calculated using the Kruskal-Wallis sum rank test (alpha = 0.05) and it detected differential abundant features at genus and species level within three sample groups. The taxonomic-level significance was

139 then tested using the pairwise Wilcoxon rank-sum tests (alpha = 0.05). Finally, the effect size 140 of each differentially abundant feature was estimated using LDA. One-against-all sample groups were compared and a linear discriminant analysis score greater than 3.6 was set as the 141 142 threshold; all-against-all sample groups were compared and a linear discriminant analysis score greater than 2.0 was set as the threshold. Cladogram was used for identification of taxa 143 at different levels of the taxonomic hierarchy between sample groups (LDA score > 2). 144 Network analysis: Network was constructed using weighted correlation network analysis or 145 weighted gene co-expression network analysis (WGCNA). Briefly, pairwise Spearman 146 147 correlation between OTUs (which was generated from LefSe analysis) was calculated using the WGCNA function. Network metrics such as betweenness, closeness, Eigen centrality, and 148 PageRank centrality of the resulting network were calculated and visualized using 'Gephi', 149 150 (https://gephi.org/) (23).

151

152 **Results** 

153 **Study design and subject attributes** 

154 The role of the microbiome in viral infections is an emerging field. We collected NP samples 155 from COVID-19 patients between 11th May 2020 and 10th October 2020 to study alterations in the NP microbiome. The schematic representation of the nasal microbiome study with 16S 156 rDNA amplicon sequencing is shown in Figure 1A. In total, 60 NP samples subjects (infected, 157 n = 47 and control, n = 13 subjects, positive and negative for SARS-CoV2 RT-qPCR test 158 159 respectively) were obtained from the Institute of Life Sciences biorepository. Out of 47 SARS-160 CoV-2-positive subjects, 25 were asymptomatic and 22 were symptomatic with mild 161 symptoms (Figure 1B). In total, 179,59,691 reads were generated. Two samples with low read

162 counts (1 from control and other from symptomatic category) were excluded and the final 163 study was performed with 58 subjects, including the control (C) [n = 12 (21%)], asymptomatic [IA, infected asymptomatic; n = 25 (43%)], and symptomatic [IS, infected symptomatic; n =164 165 21 (36%)]. The details of the participants considered for this study are described in Table 1. Differential OTUs (n = 795, p < 0.05) were obtained from a total of 3482 OTUs using the 166 deseq2 function by comparing with control NP subjects. For downstream analysis differential, 167 795 OTUs were considered. We used the t-distributed stochastic neighbor embedding (t-SNE) 168 dimension reduction method to obtain the overall distribution of NP samples with 795 OTUs 169 (Figure 1C). We found that the control and SARS-CoV-2-infected subjects showed distinct 170 segregation of OTUs in the NP microbiome, while asymptomatic and symptomatic subjects 171 showed modest separation. This indicated that the abundance of 795 differential OTUs 172 173 potentially determines the compositional distribution patterns.

#### 174 NP microbiome diversity was significantly altered in COVID-19 patients

Distinct distribution of OTUs from control and infected patients prompted us to compare the 175 evenness and richness of bacterial community compositions using Shannon and Simpson 176 alpha indices. The Shannon and Simpson alpha microbial diversity indices between control 177 178 and SARS-CoV-2-infected participants differed significantly (p-value  $\leq 0.05$ ) in pairwise Wilcoxon rank test (Shannon p-value =  $3.0 \times 10^{-4}$  and Simpson p-value =  $3.3 \times 10^{-3}$ ) (Figure 2A, 179 180 B). Although the alpha diversity indices for samples from symptomatic and asymptomatic patients compared to control subjects were found to be significantly reduced, no difference 181 182 was observed between symptomatic and asymptomatic samples (Figure 2C, D). Furthermore, 183 we used a linear regression model to establish the association between total OTU read counts for each sample and Shannon/Simpson alpha diversity indices. We found negative correlation 184

185	for both Shannon (IA - R = -0.35, $R^2 = 0.44$ , p = 0.083; IS - R = -0.54, $R^2 = 0.48$ , p = 0.012)
186	and Simpson (IA - R = -0.58, R <sup>2</sup> = 0.68, p = 0.0028; IS - R = -0.77, R <sup>2</sup> = 0.63, p = $7.7 \times 10^{-5}$ )
187	alpha diversity indices with 95% confidence intervals with total OTU counts (Figure 2E, F).
188	To further understand the microbial composition dissimilarity within the samples, we analyzed
189	beta diversity using principal coordinate analysis (PCoA) and applied both unweighted
190	(microbial richness) and weighted (microbial richness and abundance) unifrac distance
191	methods. The first two components of PCoA showed 60.3% and 80.1% variance for the
192	unweighted and weighted unifrac method. The overall difference in microbial population
193	showed two different clusters of control and SARS-CoV-2-infected patients (IA and IS) in the
194	unifrac weighted method, while the unifrac unweighted method showed more clear
195	segregation between symptomatic and asymptomatic samples (Figure 2G). We assessed the
196	significance of beta diversity to calculate unifrac distance matrix (PERMANOVA test with
197	999 permutations) for both unweighted and weighted methods and found that the three sample
198	groups (C, IA, and IS) differed significantly (P = 0.001) with 18% variance explained ( $R^2$ =
199	0.18842).

#### 200 NP microbiome dysbiosis in COVID-19 patients

Alterations in the microbial diversity prompted us to determine microbial dysbiosis index (DI) (alterations in the microbial community) across the three groups (C, IA, and IS). We performed PCoA using the Bray Curtis distance matrix and found that NP microbiota was significantly altered (p = 0.001) with 61% variation in distances explained ( $R^2 = 0.6136$ ) assessed by ADONIS test. Next, we calculated the Euclidean distance from the centroid for samples from control (median = 0.3404), asymptomatic (median = 0.1881) and symptomatic (median = 0.1511) individuals and calculated the DI (Supplementary figure 1B). The overall 208 observed DI was significant (Kruskal-Wallis test, p = 1.317E-07) across all the groups. 209 Pairwise comparison showed significant dysbiosis between control vs symptomatic (p =5.6E–09) and control vs. asymptomatic (p = 1.1E-09) groups; however, dysbiosis between 210 211 asymptomatic and symptomatic (p = 0.016) pair was not highly significant (Figure 2H). We also observed highly significant dysbiosis (p = 2.2E-12) between the control and infected 212 group (Supplemental Figure 1A, 1C). This showed that compared to that in the control 213 214 subjects, the NP microbial community is severely altered in both symptomatic and asymptomatic COVID-19 patients. 215

# Distinct microbial composition and abundance at phylum and family levels in patients suffering from SARS-CoV2 infection

The alpha and beta diversities, and DI showed that the NP microbiome was significantly 218 altered in COVID-19 patients. Next, we aimed to identify the microbial communities that were 219 altered at the phylum and family levels in three sample groups. We found 795 differential 220 221 OTUs, out of which, 12 phyla, 65 orders, 126 families, and 240 genera were present in all 222 three groups (C, IA, and IS) (Supplemental Table 1). The 12 phyla and their significance is 223 shown in Table 2. The most significant bacteria in phylum level were Actinobacteria (p =224 9.96E-07) and Proteobacteria (p = 9.61E-07), including 9 other phyla assessed using the Kruskal-Wallis test. The abundance of phyla Firmicutes (p = 4.65E-02) and Actinobacteria (p225 = 9.96E-07) were significantly higher in the SARS-CoV-2-infected groups (symptomatic and 226 asymptomatic). In contrast, Bacteroidetes (p = 1.48E-06) and Proteobacteria (p = 6.56E-07) 227 were highly abundant in the control group (non-infected) (Supplemental Figure 2A). 228 229 Furthermore, we analyzed relative abundance of top 30 families and found enrichment of 230 *Mycobacteriaceae*, *Propionibaceriaceae*, and *Streptomycetaceae* (Supplemental Figure 2B).

These families contain opportunistic pathogens in both symptomatic and asymptomatic COVID-19 patients, while these families are absent in control subjects. Top families and their significance is shown in Table 3.

# Taxonomic classifications based on OTU abundance showed sample group segregation at the genus level

236 To further our understanding regarding the 795 differentially abundant OTUs, we used the NMDS approach at phylum, order, family, and genus levels for C, IA, and IS sample groups 237 using the Bray-Curtis distance matrix. Statistical significance using ANOSIM for phylum (R 238 239 = 0.262, p = 1.7E-03), order (R = 0.322, p = 3E-04), family (R = 0.3461, p = 3E-04), and genus (R = 0.3507, p = 3E-04) showed gradual increase in R-value for genus. This indicated 240 that as we go lower in the taxonomic classification, the variance in the OTUs provides better 241 sample segregation. The differential OTUs present at the genus level in three sample groups 242 have high level of dissimilarity (35%) with R = 0.3507 and show clear sample segregation 243 244 (Figure 3A). To further validate the NMDS findings and identify the NP OTU differences between C, IA, and IS sample groups, we used sample correlation (Spearman matrix) 245 (Supplemental Figure 3A-D). The sample correlation matrix clearly showed distinction among 246 247 C, IA, and IS with respect to taxonomic classification (Figure 3B). To further reconcile the distinct sample segregation at higher to lower taxonomic level based on OTU abundances, we 248 plotted density histogram of correlation coefficient values (obtained in sample correlation). 249 250 The mean and median value of each density plot revealed lack of difference between the C, IA and IS groups at the phylum level. Furthermore, subtle differences were observed at the 251 order and family level. However, at the genus level, we found comprehensible differences 252 between C (mean = 7.95E-01; median = 6.39E-01), IA (mean = 5.65E-01; median = 8.33E-253

254 01) and IS (mean = 6.51E-01; median = 7.01E-01) (Table 4) (Figure 3B). To evaluate the 255 statistical significance of densities based on sample segregation, we calculated cumulative distribution distance (D) and significance between C, IA, and IS groups using the 256 257 Kolmogorov-Smirnov (KS) test for each taxonomic rank (Table 5). We observed that compared to that at other taxonomic levels, all the comparisons were highly significant at the 258 genus level. Based on the 'D' value comparison the samples were well distributed in C vs. IA 259 260 (D = 5.94E-01; p-value < 2.2E-16), C vs. IS (D = 5.06E-01; p = 3.308E-14), and IA vs. IS (D = 5.94E-01; p-value < 2.2E-16), C vs. IS (D = 5.06E-01; p = 3.308E-14), and IA vs. IS (D = 5.94E-01; p = 3.= 2.28E-01; p = < 2.2E-16) at the genus level. The overall sample distribution differences 261 were highly enriched at the genus level than between these three groups. Although the 'D' 262 value between IA and IS groups was less but the distribution pattern showed significant 263 differences between them. The above observance enables us to consider the genus level OTUs 264 265 (n = 240) for downstream analysis.

#### 266 Cluster-specific OTUs at genus level identified unique sample-specific OTUs

267 To gain insight regarding how the bacterial genera were segmented among three groups of samples, we performed genus level OTU correlation (n = 240) and calculated the correlation 268 coefficient (Spearman correlation), followed by unsupervised hierarchical clustering (Figure 269 270 3C). We identified 5 distinct clusters, C1 (n = 23), C2 (n = 109), C3 (n = 59), C4 (n = 33) and C5 (n = 16), with variable number of OTUs (Supplemental Table 2). The heat map 271 corresponding to each cluster is shown in Supplemental Figure 4. Cluster-wise OTUs relative 272 abundance density maps were constructed, which distinguished OTUs that were enriched in 273 274 IA/IS (C1, C3, C4, and C5) and in control samples (C2) (Figure 3D). Some of the enriched 275 cluster-specific OTUs in IA/IS are Mycobacterium (1763), Mycolicibacterium (1766), Mycobacteroides (1774), Halothiobacillus (927), Flavobacterium (986), Bifidobacterium 276

(1695), Streptomyces (1884), Rothia (2047), and Mycoplasma (2100). C2, a control-specific 277 278 cluster, contains OTUs such as Thermomicrobium (500), Kingella (502), Enterobacter (547), Bacteroides (821), and Prevotella (840). Thus, this analysis shows the distinction in genus-279 280 specific OTUs for both SARS-CoV-2-infected and control subjects. Next, we performed CCA on each of the clusters (C1 to C5) to eliminate sample heterogeneity and enhance the 281 stringency of our analysis pipeline (Figure 3E). We considered the first two components of 282 283 CCA that explained cumulative variance for the clusters. Cluster C4 explained the highest cumulative variance of 94.8%, while cluster C3 showed the lowest variance of 12.9% (Figure 284 285 3F, Supplemental Figure 5A-E). CCA showed efficient sample clustering, which is reminiscent of the density plot (Figure 3D-E). 286

#### 287 LefSe analysis identified unique OTUs at genus level in COVID-19 patients

The CCA analysis prompted us to select clusters with maximum variance explained. 288 Therefore, we considered all the clusters with  $\geq 30\%$  variance, which includes all the clusters 289 290 except C3, for LDA. OTUs (n = 181) were extracted from clusters (C1, C2, C4, C5) and plotted 291 in a heat map with their abundance (Figure 4A). Different genera could be clearly 292 distinguished between C, IA, and IS sample groups. Next, we performed LefSe to distinguish 293 the most significant microbiomes from C, IA and IS groups. In a one-against-all comparison (C with IA and IS), we got 40 genera in a control group, 34 genera in the symptomatic group, 294 and 4 genera in an asymptomatic group (LDA score  $\lceil \log 10 \rceil > 3.6$ ). The genera obtained from 295 296 one-against-all are highlighted in heatmap (Figure 4A and Supplemental Figure 6A). Relative 297 abundance of each OTU obtained from C, IA, and IS groups are shown in stack plots with clear segregation in OTUs for individual samples (Supplemental Figure 6B-D). The DI 298 299 calculated from these genera showed high dysbiosis between control and SARS-CoV-2-

medRxiv preprint doi: https://doi.org/10.1101/2021.11.10.21266147; this version posted November 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted medRxiv a license to display the preprint in perpetuity. All rights reserved. No reuse allowed without permission.

300 infected patients (Supplemental Figure 6E-F). We further increased the LefSe stringency by 301 using all-against-all (each sample group compared with each other) comparisons and constructed a cladogram and a bar plot (Figure 4B-C). All the genera obtained from LefSe 302 303 (One against all and all against all) with their LDA scores and comparison are listed in Supplemental Table 3 and Table 6. We obtained 12 significantly enriched genera of 304 Gallibacterium, Orientia, Acidocella, and Citrobacter in control samples (LDA score 305 306 [log10] > 2.0, Mycoplasma, Streptosporangium, Mycobacterium, Mycolicibacterium, Mycolicibacillus, and Mycobacteroides in symptomatic samples, and Oerskovia and 307 308 *Cellulosimicrobium* in asymptomatic samples (Figure 4C). The histogram showing the relative abundance of the 12 genera for each C, IA, and IS sample group clearly distinguishes each 309 sample type (Figure 4D). Finally, we used weighted correlation network analysis to construct 310 a network (Spearman correlation) with 12 genera identified using the LDA analysis. The 311 network creates two distinct modules, one for control groups and another for both symptomatic 312 and asymptomatic groups. We obtained strong correlation within the genera of C, IA, and IS 313 314 sample groups (Table 7). However, the correlation between C vs. IA was extremely weak and correlation was not obtained for C vs. IS groups. The network analysis suggested that the NP 315 microbiota of the control group was clearly distinct from that of the asymptomatic and 316 symptomatic groups. The DI of the 12 genera showed the highest significance between C vs. 317 IS (p = 4.7E-05), while significant dysbiosis was not observed between IA and IS groups 318 319 (Figure 4F). Overall, our analysis confirms the significance of the genera identified and their associations with symptomatic and asymptomatic COVID-19 patients. 320

#### 321 Distinct correlation of OTUs with clinical symptoms in COVID-19 patients

322 To evaluate the accuracy of LDA classification that identified eight bacterial genera in the IA and IS sample group, we tested the ROC (receiver operating characteristics) - AUC (area 323 under the curve) score. We obtained a value of 0.8 with 95% confidence interval for true 324 325 positive classification, showing 80% sensitivity and specificity of data obtained from LDA analysis (Figure 5A). Next, we used the Spearman correlation matrix to identify the 326 association of symptoms with the genera. Interestingly, chest pain showed high positive 327 328 correlation with Mycoplasma, Mycobacterium, Mycolicibacterium, Mycolicibacillus, and Mycobacteroides which were related to IS group, and weak correlation with Oerskovia and 329 330 *Cellulosimicrobium*, which were associated with the IA group. *Mycoplasma*, however, showed a strong correlation with both chest pain (0.4446) and fever (0.4214) (Figure 5B). 331 ROC-AUC analysis for chest pain and fever showed 0.90 and 0.79 scores, respectively, with 332 333 eight bacterial genera (Figure 5C-D). We extended our study at the species level for the 12 genera found in LDA analysis and observed that 54 species were represented in a heat map 334 for C, IA, and IS group (Supplemental Figure 7A). Several known opportunistic pathogens 335 336 such as Mycobacterium tuberculosis, Mycobacterium avium, and Mycoplasma pneumonia were highly abundant in the SARS-CoV-2-infected patients. The significance of the 54 337 338 bacterial species was assessed using the Kruskal-Wallis test and the top 30 significant species were plotted in the bubble plot (Supplemental Figure 7B). In sum, we established the 339 association of pathogenic microbes with COVID-19 disease and showed susceptibility to 340 alterations in the NP microbiome in case of infection in SARS-COV-2. We also identified the 341 compositional difference in NP microbiota between symptomatic and asymptomatic group. 342

343 **Discussion** 

344 Scientists worldwide are trying to understand the pathophysiology of SARS-CoV-2 infection and the associated alterations in the host, including those in the microbiome. As SARS-CoV-345 2 infection initiates in the upper respiratory tract, we investigated the alterations in the NP 346 347 microbiota of COVID-19 patients. We amplified the 16rRNA gene of variable regions (V1-V9) and performed long-read sequencing using Oxford Nanopore technology. Subsequently, 348 we have used multiple bioinformatics approaches to cross-validate our data sets at various 349 350 levels and identify the most significant bacterial population in the NP microbiome of COVID-19 patients. We found significant changes in abundance, diversity, and DI of SARS-CoV-2-351 352 infected patients compared to those of the control. The IA and IS groups also showed overall significant alterations in microbiota composition. We found abundance of opportunistic 353 pathogens such as Mycoplasma and Mycobacterium in symptomatic patients, which correlated 354 355 strongly with patient symptoms such as chest pain and fever. Insights into species level abundance revealed the presence of Mycoplasma pneumoniae, Mycobacterium tuberculosis, 356 Mycobacterium avium, and Mycolicibacterium sp. in the SARS-CoV-2-infected patients. To 357 358 the best of our knowledge, this is the first comprehensive study to report abundance of 359 opportunistic pathogens such as Mycoplasma pneumoniae and Mycobacterium tuberculosis 360 based on the complete sequence of the 16S rRNA variable regions in patients with SARS-CoV-2 infection. 361

Respiratory infections alter the NP microbiota, which reduces the diversity of the NP microbial ecosystem and promotes the growth of opportunistic pathogens (24). At the phylum level, Proteobacteria, Firmicutes, and Actinobacteria were detected in all NP samples. However, the abundance of Firmicutes and Actinobacteria was significantly higher in both symptomatic and asymptomatic groups. Our results are in partial agreement with those of Ventero et al., who

367 found the abundance of Firmicutes, Bacteroidota, Proteobacteria, and Actinobacteria in the NP samples of COVID-19 patients (13). Only few studies have shown either no alterations or 368 significant changes in the microbiome composition of the nasopharynx during COVID-19 369 370 infection. Maio et al. and Braun et al. did not find any significant alterations in NP microbial composition (12, 25). However, other studies showed prevalence of opportunistic pathogens 371 such as Staphylococcus, Anelloviridae, Pseudomonas, Haemophilus, Stenotrophomonas, 372 Redondoviridae, and Pseudomonas aeruginosa in COVID-19 patients (11, 13, 15-18). 373 Compared to earlier reports, our study also revealed overall changes in the composition of the 374 375 NP microbial community, reduction in bacterial diversity due to COVID-19 infection and the 376 presence of opportunistic pathogens such as *Mycoplasma* and *Mycobacterium* in COVID-19 patient cohort. 377

378 Most of the NP microbial studies amplify short 16S rRNA gene using the Illumina platform, which is more accurate but is limited by taxonomic resolution owing to sequencing of shorter 379 reads (26) and sequencing of the specific variable region. The taxonomic resolution can be 380 improved to genus, species, and even at the strain level by sequencing the V1–V9 (~ 1600 bp) 381 variable regions of the 16S rRNA gene (26, 27). In this study, we have used the Oxford 382 Nanopore<sup>™</sup> long read sequenced platform and sequenced V1 to V9 (~1.6 kb) of 16S variable 383 regions and successfully obtained taxonomic resolution to genus and some extent species 384 level. This has provided us immense advantage of determining the abundance of opportunistic 385 pathogens in the NP of the COVID-19 patients. Until now, only Mostafa et al. has used 386 metagenomics for COVID-19 NP samples using Oxford Nanopore technology. They have 387 sequenced both RNA and DNA from the NP samples without any PCR amplification. They 388 389 not only identified the SARS-CoV-2 virus in the samples but also potential pathogens that 390 may lead to co-infections (18). Our study is the first to use 16S amplification of  $\sim$ 1.6 Kb 391 variable regions to identify the bacterial community associated with infected and control NP. However, 16S rRNA gene amplification may introduce PCR biases, however, more subjects 392 393 and a robust analysis pipeline may dilute these biases. This study is the first comprehensive study from Odisha cohort and second from India. Gupta et al. used the Illumina platform for 394 16S amplicon sequencing and found enrichment of several opportunistic pathogens (17). 395 396 Interestingly, Mycoplasma, Mycolicocibacterium, and Mycobacterium were not present in their list. This could be due to the analysis pipeline or region-specific differences. 397 398 Nevertheless, the identification of opportunistic pathogens and their increase in abundance in COVID-19 patients is one of the important aspects of this study. 399

Our comprehensive bioinformatics analysis with sample and OTU correlation analysis 400 401 distinguished COVID-19 infected and control samples at the genera level. Furthermore, LDA analysis identified a significantly high abundance of Mycobacterium and Mycoplasma in 402 symptomatic patients, which correlated well with the occurrence of fever and chest pain. 403 404 Significantly high relative abundance of members of family Mycobacteriaceae in the symptomatic COVID-19 group indicates the presence of both pathogenic and non-pathogenic 405 406 bacteria. Further dissection into genus level revealed the presence of several key genera, namely, Mycobacterium, Mycolicibacterium, Mycolicibacillus, and Mycobacteroides. 407 *Mycobacterium* genera are well associated with several pulmonary diseases, for example, 408 409 Mycobacterium tuberculosis is responsible for tuberculosis in humans and is associated with 410 pulmonary infection (28), while *Mycobacterium avium* is highly associated with lung disease (29). Mycolicibacterium and Mycolicibacillus are generally non-pathogenic but some species 411 have been associated with pathogenicity in humans and were isolated from hospitalized 412

patients (30). Mycobacteroides are potentially associated with soft tissue infections and 413 Mycobacteroides abscessus is a known pulmonary pathogen (30-32). Members of genus 414 Mycoplasma is a well-recognized pathogen, Mycoplasma pneumoniae being responsible for 415 416 pneumonia and other respiratory infections in humans (33, 34). Cellulosimicrobium and Oerskovia were detected in asymptomatic COVID-19 patients. Few species of 417 418 *Cellulosimicrobium* are pathogenic, although their pathogenicity was not clear under normal conditions because those species were isolated from hospitalized patients with acute renal 419 failure (35, 36). Results of our and other reports have proved the association of opportunistic 420 421 pathogens with alterations in the diversity of the microbial communities in symptomatic and asymptomatic COVID-19 patients. This study establishes a new set of opportunistic pathogens 422 in the context of NP microbiome in COVID-19 infected patients. Moreover, this study clearly 423 424 distinguishes between the NP microbial composition of symptomatic and asymptomatic groups using LefSe with AUC-ROC validation. Thus, we believe that SARS-CoV-2 virulence 425 may promote the growth of opportunistic pathogens and may lead to coinfection or secondary 426 427 infection in COVID-19 patients.

Our study has certain limitations. The subject size is limited and a larger cohort would have
strengthened our findings. The clinical manifestations are limited, and therefore, the larger
picture is difficult to interpret. Future studies should include NP samples of vaccinated,
asymptomatic, and hospitalized COVID-19 patients with detailed pathophysiology.
Furthermore, blood biochemistry and metabolite studies from the serum would boost
conclusions regarding functional aspects of the NP microbiome.

434 Author Contribution

435	P.P. and A.P. conceptualized the study and secured funding. P.P and S.M. initiated the work,
436	directed overall workflow, interpreted data, and troubleshoot the experiments. R.M. did most
437	of the bioinformatics analysis and S.M., K.C.M. and A.G. helped in bioinformatics analysis
438	and troubleshooting. S.M., S.A., M.S., P.M., R.Y., H.D., S.A.A., S.D., and D.S., helped with
439	the preprocessing of the samples in Biosafety level 3 (BSL3) facility and nucleic acid
440	extractions. S.M., S.A., M.S., P.M., and R.Y. were involved in amplicon library preparations.
441	R.S. provided samples from the Biorepository. A.S., R.D., S.S., T.K.B., S.C., G.H.S., R.S.,
442	S.K.R., P.P., and A.P. coordinated with COVID-19 sampling and testing at BSL3. P.P., S.M.,
443	and R.M wrote the manuscript.
444	Conflict of interest
445	The authors declare no competing commercial or financial interests in relation to this work.
446	Acknowledgments
446 447	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT),
446 447 448	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project
446 447 448 449	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re-
446 447 448 449 450	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge
446 447 448 449 450 451	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge Biorepository, BSL-3, and BSL-2 laboratories, qPCR, and DNA-sequencing institutional
446 447 448 449 450 451 452	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge Biorepository, BSL-3, and BSL-2 laboratories, qPCR, and DNA-sequencing institutional central core facilities. R.M., S.M., and K.C.M received their fellowships from
446 447 448 449 450 451 452 453	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge Biorepository, BSL-3, and BSL-2 laboratories, qPCR, and DNA-sequencing institutional central core facilities. R.M., S.M., and K.C.M received their fellowships from Ramalingaswami, ILS Flagship, and SERB core research grant, respectively. We thank all the
446 447 448 449 450 451 452 453 454	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge Biorepository, BSL-3, and BSL-2 laboratories, qPCR, and DNA-sequencing institutional central core facilities. R.M., S.M., and K.C.M received their fellowships from Ramalingaswami, ILS Flagship, and SERB core research grant, respectively. We thank all the volunteers who provided samples for research purposes.
446 447 448 449 450 451 452 453 454	Acknowledgments We acknowledge the institute's core funding from the Department of Biotechnology (DBT), Government of India. This work was also supported by the ILS flagship project (BT/ILS/Flagship/2019) from DBT, Ramalingaswami Re-entry fellowship (BT/RLF/Re- entry/25/2015), and SERB core research grant (CRG/2018/002052). We also acknowledge Biorepository, BSL-3, and BSL-2 laboratories, qPCR, and DNA-sequencing institutional central core facilities. R.M., S.M., and K.C.M received their fellowships from Ramalingaswami, ILS Flagship, and SERB core research grant, respectively. We thank all the volunteers who provided samples for research purposes.

- He Y, Wang J, Li F, Shi Y. Main Clinical Features of COVID-19 and Potential
  Prognostic and Therapeutic Value of the Microbiota in SARS-CoV-2 Infections. Front
  Microbiol. 2020;11:1302.
- 462 2. Khatiwada S, Subedi A. Lung microbiome and coronavirus disease 2019 (COVID-19):
- 463 Possible link and implications. Hum Microb J. 2020;17:100073.
- 3. Zou X, Chen K, Zou J, Han P, Hao J, Han Z. Single-cell RNA-seq data analysis on the
- receptor ACE2 expression reveals the potential risk of different human organs vulnerable to
- 466 2019-nCoV infection. Front Med. 2020;14(2):185-92.
- 467 4. Hou YJ, Okuda K, Edwards CE, Martinez DR, Asakura T, Dinnon KH, 3rd, et al. SARS-
- 468 CoV-2 Reverse Genetics Reveals a Variable Infection Gradient in the Respiratory Tract. Cell.
  469 2020;182(2):429-46 e14.
- 470 5. Mlcochova P, Kemp SA, Dhar MS, Papa G, Meng B, Ferreira I, et al. SARS-CoV-2
- 471 B.1.617.2 Delta variant replication and immune evasion. Nature. 2021.
- 472 6. Belkaid Y, Harrison OJ. Homeostatic Immunity and the Microbiota. Immunity.
  473 2017;46(4):562-76.
- Fazlollahi M, Lee TD, Andrade J, Oguntuyo K, Chun Y, Grishina G, et al. The nasal
  microbiome in asthma. J Allergy Clin Immunol. 2018;142(3):834-43 e2.
- 476 8. de Steenhuijsen Piters WA, Sanders EA, Bogaert D. The role of the local microbial
  477 ecosystem in respiratory health and disease. Philos Trans R Soc Lond B Biol Sci.
  478 2015;370(1675).
- 479 9. Kumpitsch C, Koskinen K, Schopf V, Moissl-Eichinger C. The microbiome of the upper
  480 respiratory tract in health and disease. BMC Biol. 2019;17(1):87.

481	10. Yildiz S, Mazel-Sanchez B, Kandasamy M, Manicassamy B, Schmolke M. Influenza A
482	virus infection impacts systemic microbiota dynamics and causes quantitative enteric
483	dysbiosis. Microbiome. 2018;6(1):9.
484	11. Engen PA, Naqib A, Jennings C, Green SJ, Landay A, Keshavarzian A, et al.
485	Nasopharyngeal Microbiota in SARS-CoV-2 Positive and Negative Patients. Biol Proced
486	Online. 2021;23(1):10.
487	12. De Maio F, Posteraro B, Ponziani FR, Cattani P, Gasbarrini A, Sanguinetti M.
488	Nasopharyngeal Microbiota Profiling of SARS-CoV-2 Infected Patients. Biol Proced Online.
489	2020;22:18.
490	13. Ventero MP, Cuadrat RRC, Vidal I, Andrade BGN, Molina-Pardines C, Haro-Moreno
491	JM, et al. Nasopharyngeal Microbial Communities of Patients Infected With SARS-CoV-2
492	That Developed COVID-19. Front Microbiol. 2021;12:637430.
493	14. Feehan AK, Rose R, Nolan DJ, Spitz AM, Graubics K, Colwell RR, et al.
494	Nasopharyngeal Microbiome Community Composition and Structure Is Associated with
495	Severity of COVID-19 Disease and Breathing Treatment. Applied Microbiology.
496	2021;1(2):177-88.
497	15. Rhoades NS, Pinski AN, Monsibais AN, Jankeel A, Doratt BM, Cinco IR, et al. Acute
498	SARS-CoV-2 infection is associated with an increased abundance of bacterial pathogens,
499	including Pseudomonas aeruginosa in the nose. Cell Rep. 2021;36(9):109637.
500	16. Nardelli C, Gentile I, Setaro M, Di Domenico C, Pinchera B, Buonomo AR, et al.
501	Nasopharyngeal Microbiome Signature in COVID-19 Positive Patients: Can We Definitively
502	Get a Role to Fusobacterium periodonticum? Front Cell Infect Microbiol. 2021;11:625581.

503	17. Gupta A, Karyakarte R, Joshi S, Das R, Jani K, Shouche Y, et al. Nasopharyngeal
504	microbiome reveals the prevalence of opportunistic pathogens in SARS-CoV-2 infected
505	individuals and their association with host types. Microbes Infect. 2021:104880.

- 18. Mostafa HH, Fissel JA, Fanelli B, Bergman Y, Gniazdowski V, Dadlani M, et al.
- 507 Metagenomic Next-Generation Sequencing of Nasopharyngeal Specimens Collected from
- 508 Confirmed and Suspect COVID-19 Patients. mBio. 2020;11(6):1-13.
- 509 19. Liu J, Liu S, Zhang Z, Lee X, Wu W, Huang Z, et al. Association between the
- 510 nasopharyngeal microbiome and metabolome in patients with COVID-19. Synth Syst
- 511 Biotechnol. 2021;6(3):135-43.
- 512 20. Haiminen N, Utro F, Seabolt E, Parida L. Functional profiling of COVID-19 respiratory
  513 tract microbiomes. Sci Rep. 2021;11(1):6433.
- 514 21. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2. Genome
  515 Biol. 2019;20(1):257.
- 516 22. Lloyd-Price J, Arze C, Ananthakrishnan AN, Schirmer M, Avila-Pacheco J, Poon TW,
- et al. Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. Nature.
  2019;569(7758):655-62.
- 519 23. Bastian M, Heymann S, Jacomy M. Gephi: An Open Source Software for Exploring and
  520 Manipulating Networks. Proceedings of the International AAAI Conference on Web and
  521 Social Media. 2009;3(1):361-2.
- 522 24. Santacroce L, Charitos IA, Ballini A, Inchingolo F, Luperto P, De Nitto E, et al. The
- 523 Human Respiratory System and its Microbiome at a Glimpse. Biology (Basel). 2020;9(10).

524	25. Braun T, Halevi S, Hadar R, Efroni G, Glick Saar E, Keller N, et al. SARS-CoV-2 does
525	not have a strong effect on the nasopharyngeal microbial composition. Sci Rep.
526	2021;11(1):8922.
527	26. Johnson JS, Spakowicz DJ, Hong BY, Petersen LM, Demkowicz P, Chen L, et al.
528	Evaluation of 16S rRNA gene sequencing for species and strain-level microbiome analysis.
529	Nat Commun. 2019;10(1):5029.
530	27. Kaul D, Rathnasinghe R, Ferres M, Tan GS, Barrera A, Pickett BE, et al. Microbiome
531	disturbance and resilience dynamics of the upper respiratory tract during influenza A virus
532	infection. Nat Commun. 2020;11(1):2537.
533	28. Peto HM, Pratt RH, Harrington TA, LoBue PA, Armstrong LR. Epidemiology of
534	extrapulmonary tuberculosis in the United States, 1993-2006. Clin Infect Dis.
535	2009;49(9):1350-7.
536	29. Hwang JA, Kim S, Jo KW, Shim TS. Natural history of Mycobacterium avium complex
537	lung disease in untreated patients with stable course. Eur Respir J. 2017;49(3):1600537.
538	30. Gupta RS, Lo B, Son J. Phylogenomics and Comparative Genomic Studies Robustly
539	Support Division of the Genus Mycobacterium into an Emended Genus Mycobacterium and
540	Four Novel Genera. Front Microbiol. 2018;9:67.
541	31. Batchelder HR, Story-Roller E, Lloyd EP, Kaushik A, Bigelow KM, Maggioncalda EC,
542	et al. Development of a penem antibiotic against Mycobacteroides abscessus. Commun Biol.
543	2020;3(1):741.
544	32. Tortoli E. Microbiological features and clinical relevance of new species of the genus
545	Mycobacterium. Clin Microbiol Rev. 2014;27(4):727-52.

546	33. Beeton ML, Zhang XS, Uldum SA, Bebear C, Dumke R, Gullsby K, et al. Mycoplasma
547	pneumoniae infections, 11 countries in Europe and Israel, 2011 to 2016. Euro Surveill.
548	2020;25(2).
549	34. Foy HM. Infections caused by Mycoplasma pneumoniae and possible carrier state in
550	different populations of patients. Clin Infect Dis. 1993;17 Suppl 1:S37-46.
551	35. Sharma A, Gilbert JA, Lal R. (Meta)genomic insights into the pathogenome of
552	Cellulosimicrobium cellulans. Sci Rep. 2016;6:25527.
553	36. Delport J, Wakabayashi AT, Anantha RV, Lannigan R, John M, McCormick JK.
554	Cellulosmicrobium cellulans isolated from a patient with acute renal failure. JMM Case
555	Reports. 2014;1(2):e000976.
556	Figure Legends:
557	Figure 1: Schema of nasopharyngeal sample processing, 16S sequencing, and OTU-
558	<b>based sample distribution.</b> (A) Flow chart showing nasopharyngeal sample processing for
559	DNA extraction, amplicon library preparation, Oxford Nanopore <sup>™</sup> sequencing, and
560	bioinformatics analysis pipeline. (B) Pie chart showing nasopharyngeal samples (controls,
561	symptomatic, and asymptomatic) used in this study. (C) t-SNE plot showing the OTU-based
562	sample distribution and ordination points for control, symptomatic, and asymptomatic
563	samples.
564	
565	Figure 2: Alpha/beta diversities and dysbiosis index in COVID-19-positive and negative
566	nasopharyngeal sample. (A-B) Alpha diversity index (Shannon/Simpson) between control
567	and COVID-19-infected samples (pairwise Wilcoxon rank-sum test $p = \leq 0.05$ ). (C-D) Same

analysis as above where the COVID-19-infected samples are classified as asymptomatic and

569 symptomatic compared to the control group. (E-F) Linear regression model showing the 570 association between total OTU count and Shannon/Simpson diversity index for each sample; the shaded grey region represents 95% confidence intervals of two groups, symptomatic and 571 572 asymptomatic, with correlation (Spearman) regression line [Shannon: R = -0.35(asymptomatic), R = -0.54 (symptomatic) and Simpson: R = -0.58 (asymptomatic), R = -0.77573 (symptomatic)]. (G) Principal coordinate analysis (PCoA) showing beta diversity in 574 symptomatic, and control sample unifrac 575 asymptomatic, groups based on (weighted/unweighted) distance (p = 0.001, PERMANOVA). (H) Violin plot showing 576 dysbiosis indexes of samples from control, asymptomatic, and symptomatic participants 577 578 (pairwise Wilcoxon rank-sum test  $p = \le 0.05$ ).

579

580 Figure 3: Taxonomic classification of bacterial communities using non-metric multidimensional scaling (NMDS), correlation, and canonical correspondence analysis 581 582 (CCA). (A) NMDS ordination of Bray-Curtis distance matrix based on all samples and 583 bacterial communities of each taxonomy level (phylum, order, family, and genus) (ANOSIM  $p = \langle 0.05 \rangle$ . (B) The density plot representing the Spearman correlation coefficient at each 584 585 taxonomy level (phylum, order, family, and genus); dotted line indicates the mean value of each sample group (Kolmogorov–Smirnov (KS) Test  $p \le 0.05$ ). (C) Heat map of Spearman 586 correlation for genus level with sample correlation (lower) and OTU correlation (upper). Five 587 588 clusters (C1, C2, C3, C4, and C5) were generated using unsupervised hierarchical clustering 589 from the OTU correlation plot. (D) Sample-wise OTU density plot for each cluster (C1, C2, C3, C4, and C5) showing relative abundance. (E-F) CCA plot of microbial community 590 591 composition for each cluster and bar plot representing cumulative variation percentage from

two components [C1 (92.4%), C2 (80.5%), C3 (12.9%), C4 (94.8%), and C5 (39.4%)]. Dotted
line shows 30% variance cut-off for downstream analysis.

594

595 Figure 4: Linear discriminant analysis effect size (LefSe) analysis reveals distinct genuslevel OTUs in control, asymptomatic and symptomatic. (A) Heat map showing genus level 596 OTU (n = 181) abundance distribution from four clusters (C1, C2, C4, and C5) identified from 597 598 CCA analysis in control, asymptomatic and symptomatic samples. The OTUs marked on either side of the heat map were obtained from one-against-all and all-against-all comparison 599 600 in LDA analysis (B) The cladogram shows the output of the LEfSe (LDA score >2.0), which identifies taxonomic differences between sample groups. Each circle represents a bacterial 601 taxon, and each ring of taxonomy level starting with kingdom in the innermost circle is 602 603 followed by phylum, class, order, family, and genus in the outermost circle. The different color intensities indicate the different taxonomy levels and the diameter of each circle is 604 proportional to the taxon's abundance and correlates with the LDA score. (C) The histogram 605 606 of the LDA scores (score >2.0 and all-against-all) was computed for differentially abundant taxa between sample groups. The effect size of specific taxa in the particular group at the 607 608 genus level. (D) Histogram of the all LefSe-specific taxa (Mycoplasma, Streptosporangium, *Mycolicibacterium*, 609 Citrobacter, Acidocella, *Mycolicibacillus*, *Mycobacterium*, Mycobacteroides, Orientia, Gallibacterium, Cellulosimicrobium, and Oerskovia) showing 610 611 relative abundance across sample groups. Solid and dotted lines show median and mean 612 relative abundance respectively. (E) Weighted correlation network analysis (WGCNA) was used for network construction and plotted using Gephi. Each node of the network represents 613 614 the individual bacterial genera with their respective abundance size and the edges represent

615	correlation strength with edge weight by thickness. The pie chart within each node represents
616	abundance for each genus. The dotted line shows two distinct modules (control and infected)
617	created in the network analysis. (F) Violin plot showing the dysbiosis indexes of LefSe sample
618	groups (pairwise Wilcoxon rank-sum test p-value $< 0.05$ ).
619	
620	Figure 5: Area under the curve-receiver operating characteristic (AUC-ROC) validation
621	and correlation of genera with the symptoms of COVID-19 subjects. (A) ROC curve for
622	LDA classified symptomatic and asymptomatic group. The AUC w 0.80 with a 95%
623	confidence interval (CI). (B) Correlation between bacteria at genus level and clinical
624	symptoms of patients. (C-D) ROC curve for chest pain and fever in the symptomatic and
625	asymptomatic group. The AUCs were 0.904 (chest pain) and 0.793 (fever) with a 95%
626	confidence interval (CI).
627	
628	
629	
630	
631	
632	
633	
634	
635	
636	
637	

	Control (n-12)	Asymptomatic	Symptomatic (n=21)
Sex	Control (II=12)	(n-23)	(11-21)
Male	5 (41.66%)	18 (72%)	19 (90.47%)
Female	7 (58.33%)	7 (28%)	2 (9.52%)
Age (years)	31 (median)	26 (median)	32 (median)
Symptoms			
Dry Cough	NA	NA	7 (33.3%)
Fever	NA	NA	17 (80.95%)
Tiredness	NA	NA	8 (38.09%)
Sore throat	NA	NA	11(52.38%)
Body pain	NA	NA	13(61.9%)
Chest pain	NA	NA	9(42.85%)
Fever + Body pain	NA	NA	4(19.04%)
Fever + multiple symptoms*	NA	NA	16(76.19%)
Fever + chest pain	NA	NA	8(38.09%)
Loss of smell/taste + multiple symptoms* without fever	NA	NA	2(9.52%)

# 638Table 1: Details of samples included in this study

639

\*Multiple symptoms refer to having more than one symptoms from symptoms list.

#### **Table 2: Phylum based on relative abundance and their respective values**

		1st		3rd		
Phylum	Mean	Quartile	Median	Quartile	p_value	BH_FDR
Proteobacteria	1.51E-01	3.79E-02	4.53E-02	7.29E-02	6.56E-07	6.56E-07
Fusobacteria	2.25E-03	1.31E-03	1.72E-03	2.71E-03	9.00E-07	9.00E-07
Actinobacteria	6.83E-01	7.09E-01	7.62E-01	7.99E-01	9.96E-07	9.96E-07
Bacteroidetes	8.77E-03	5.71E-03	6.66E-03	1.00E-02	1.48E-06	1.48E-06
Tenericutes	9.37E-04	1.35E-04	5.60E-04	1.21E-03	2.21E-06	2.21E-06
Chloroflexi	2.25E-03	6.47E-04	9.84E-04	1.46E-03	6.54E-06	6.54E-06
Chlamydiae	5.93E-04	2.82E-04	4.98E-04	7.76E-04	7.27E-06	7.27E-06
Fibrobacteres	5.65E-04	2.86E-04	4.05E-04	7.10E-04	4.77E-05	4.77E-05
Thermodesulfobacteria	2.62E-02	2.09E-02	2.84E-02	3.24E-02	8.13E-04	8.13E-04
Aquificae	2.31E-04	1.19E-04	2.30E-04	2.93E-04	1.59E-03	1.59E-03
Firmicutes	1.23E-01	1.19E-01	1.28E-01	1.44E-01	4.65E-02	4.65E-02
Chlorobi	1.02E-03	6.81E-04	9.32E-04	1.27E-03	8.58E-01	8.58E-01

641

642

		1st		3rd		
Family	Mean	quartile	Median	quartile	p_value	BH_FDR
Acetobacteraceae	5.08E-02	5.00E-02	5.12E-02	5.28E-02	9.14E-07	9.14E-07
Actinomycetaceae	2.23E-02	2.10E-02	2.18E-02	2.32E-02	8.18E-03	8.18E-03
Aeromonadaceae	2.96E-02	2.54E-02	3.07E-02	3.38E-02	8.18E-03	8.18E-03
Alcaligenaceae	8.14E-02	8.50E-02	8.71E-02	9.43E-02	1.41E-06	1.41E-06
Bifidobacteriaceae	7.69E-02	4.59E-02	6.52E-02	1.04E-01	2.99E-03	2.99E-03
Brevibacteriaceae	4.37E-02	2.72E-02	3.82E-02	5.98E-02	6.35E-04	6.35E-04
Cellulomonadaceae	6.75E-02	4.60E-02	6.29E-02	8.52E-02	8.64E-05	8.64E-05
Chromobacteriaceae	2.30E-02	2.27E-02	2.30E-02	2.32E-02	9.37E-07	9.37E-07
Corynebacteriaceae	5.49E-02	4.18E-02	5.05E-02	6.52E-02	2.36E-05	2.36E-05
Enterobacteriaceae	1.26E-01	3.61E-02	1.65E-01	2.17E-01	1.73E-06	1.73E-06
Erwiniaceae	3.19E-02	2.97E-02	3.10E-02	3.46E-02	2.01E-06	2.01E-06
Erysipelotrichaceae	2.26E-02	2.21E-02	2.26E-02	2.30E-02	1.69E-06	1.69E-06
Eubacteriaceae	2.99E-02	2.58E-02	2.84E-02	3.28E-02	1.85E-06	1.85E-06
Lactobacillaceae	2.12E-02	2.05E-02	2.10E-02	2.18E-02	5.39E-03	5.39E-03
Micrococcaceae	3.52E-02	2.65E-02	3.37E-02	4.36E-02	3.49E-03	3.49E-03
Micromonosporaceae	2.41E-02	2.25E-02	2.43E-02	2.53E-02	2.65E-06	2.65E-06
Morganellaceae	5.01E-02	4.51E-02	5.30E-02	5.93E-02	1.90E-05	1.90E-05
Mycobacteriaceae	2.19E-01	1.85E-01	2.42E-01	2.69E-01	7.93E-07	7.93E-07
Neisseriaceae	5.24E-02	5.08E-02	5.23E-02	5.43E-02	1.29E-06	1.29E-06
Nocardiaceae	3.15E-02	2.46E-02	3.41E-02	3.63E-02	8.90E-07	8.90E-07
Pasteurellaceae	4.73E-02	3.34E-02	5.14E-02	5.75E-02	5.89E-05	5.89E-05
Pectobacteriaceae	2.17E-02	2.05E-02	2.15E-02	2.27E-02	2.09E-06	2.09E-06
Peptostreptococcaceae	3.01E-02	2.56E-02	2.82E-02	3.23E-02	1.89E-06	1.89E-06
Propionibacteriaceae	1.17E-01	9.00E-02	1.20E-01	1.33E-01	8.90E-07	8.90E-07
Pseudonocardiaceae	2.94E-02	2.70E-02	2.98E-02	3.27E-02	1.44E-06	1.44E-06

# 644 Table 3: Top 30 family based on relative abundance and their respective values

645

#### **Table 4: Mean and median value of density plots.**

Taxonomic				
rank	Statistic value	Control	Asymptomatic	Symptomatic
	Mean	8.96E-01	8.68E-01	8.90E-01
Phylum	Median	9.30E-01	8.95E-01	9.16E-01
	Mean	8.58E-01	7.70E-01	8.36E-01
Order	Median	8.78E-01	8.21E-01	8.93E-01
	Mean	8.22E-01	6.86E-01	7.60E-01
Family	Median	8.63E-01	7.46E-01	8.06E-01
	Mean	7.95E-01	5.65E-01	6.51E-01
Genus	Median	6.39E-01	8.33E-01	7.01E-01

# 647 Table 5: Result of Kolmogorov–Smirnov (KS) test between the densities of each

Taxonomic	Control vs	Control vs	Asymptomatic vs
rank	Asymptomatic	Symptomatic	Symptomatic
	D = 1.88e-01	D = 1.07e-01	D = 1.01e-01
Phylum	p-value = 3.03e-02	p-value = 0.4834	p-value = 9.78e-04
	D = 3.18e-01	D = 9.97e-02	D = 2.30e-01
Order	p-value = 1.237e-05	p-value = 0.5706	p-value < 2.2e-16
	D = 4.14e-01	D = 2.93e-01	D = 2.25e-01
Family	p-value = 2.706e-09	p-value = 4.75e-05	p-value < 2.2e-16
	D = 5.94e-01	D = 5.06e-01	D = 2.28e-01
Genus	p-value < 2.2e-16	p-value = 3.308e-14	p-value < 2.2e-16

# 648 taxonomic rank.

### **Table 6: Linear discriminate analysis (LDA) score for all-against-all analysis**

Genus	highest mean	Samples	LDA score	pvalue	
	among all		(log 10)		
	the classes				
Oerskovia	3.83	Asymptomatic	3.44	1.96E-02	
Cellulosimicrobium	3.85	Asymptomatic	3.5	1.73E-02	
Gallibacterium	3.84	Control	3.36	8.46E-04	
Orientia	3.88	Control	3.46	3.07E-04	
Acidocella	4.07	Control	3.81	4.68E-07	
Citrobacter	4.08	Control	3.89	2.82E-07	
Mycobacteroides	3.88	Symptomatic	3.61	3.19E-07	
Mycolicibacillus	3.98	Symptomatic	3.67	3.19E-07	
Mycolicibacterium	3.93	Symptomatic	3.71	2.55E-07	
Mycobacterium	3.9	Symptomatic	3.79	2.55E-07	
Streptosporangium	4.16	Symptomatic	3.83	1.65E-04	
Mycoplasma	4.26	Symptomatic	3.95	1.64E-06	

# **Table 7: WGCNA network data table**

B	Genus	OTUS	Module	absolute abundance	Degree	Page ranks	Eigen centrality	Modularity class	Clustering	Triangles
1	Acidocella	525	M1	1362	5	0.0673	0.4963	0	1	10
2	Citrobacter	546	M1	3463	5	0.0673	0.4963	0	1	10
3	Gallibacterium	750	M1	285	5	0.0673	0.4963	0	1	10
4	Orientia	784	M1	263	5	0.0673	0.4963	0	1	10
5	Cellulosimicrobium	1710	M2	11832	10	0.1202	1	1	0.5555	25
6	Oerskovia	1713	M2	13210	10	0.1202	1	1	0.5555	25
7	Mycobacterium	1763	M2	139017	7	0.0852	0.8208	1	0.9047	19
8	Mycolicibacterium	1766	M2	81430	7	0.0852	0.8208	1	0.9047	19
9	Mycobacteroides	1774	M2	7172	7	0.0852	0.8208	1	0.9047	19
10	Mycolicibacillus	1798	M2	2749	7	0.0852	0.8208	1	0.9047	19
11	Streptosporangium	2002	M2	163	5	0.0642	0.5732	1	1	10
12	Mycoplasma	2100	M2	475	7	0.0852	0.8208	1	0.9047	19

Figure 1









