

1 **Total Infectomes Characterization of Respiratory Infections in pre-**

2 **COVID-19 Wuhan, China**

3 **Mang Shi<sup>1,2+</sup>, Su Zhao<sup>3+</sup>, Bin Yu<sup>4+</sup>, Wei-Chen Wu<sup>1+</sup>, Yi Hu<sup>3+</sup>, Jun-Hua Tian<sup>4</sup>, Wen Yin<sup>3</sup>,**

4 **Fang Ni<sup>3</sup>, Hong-Ling Hu<sup>3</sup>, Shuang Geng<sup>3</sup>, Li Tan<sup>3</sup>, Ying Peng<sup>4</sup>, Zhi-Gang Song<sup>1</sup>, Wen**

5 **Wang<sup>1,5</sup>, Yan-Mei Chen<sup>1</sup>, Edward C. Holmes<sup>1,2</sup>, Yong-Zhen Zhang<sup>1\*</sup>**

6

7 <sup>1</sup>Shanghai Public Health Clinical Center, State Key Laboratory of Genetic Engineering, School

8 of Life Sciences and Human Phenome Institute, Fudan University, Shanghai, China.

9 <sup>2</sup>Marie Bashir Institute for Infectious Diseases and Biosecurity, School of Life and

10 Environmental Sciences and School of Medical Sciences, The University of Sydney, Sydney,

11 Australia

12 <sup>3</sup>Department of Pulmonary and Critical Care Medicine, The Central Hospital of Wuhan, Tongji

13 Medical College, Huazhong University of Science and Technology, Wuhan, China.

14 <sup>4</sup>Wuhan Center for Disease Control and Prevention, Wuhan, China.

15 <sup>5</sup>Department of Zoonosis, National Institute for Communicable Disease Control and Prevention,

16 China Center for Disease Control and Prevention, Beijing, China.

17

18 + Contributed equally to this article.

19 \* Correspondence: [zhangyongzhen@fudan.edu.cn](mailto:zhangyongzhen@fudan.edu.cn)

20 **Abstract**

21 At the end of 2019 Wuhan witnessed an outbreak of “atypical pneumonia” that later developed  
22 into a global pandemic. Metagenomic sequencing rapidly revealed the causative agent of this  
23 outbreak to be a novel coronavirus - SARS-CoV-2. Herein, to provide a snapshot of the  
24 pathogens in pneumonia-associated respiratory samples from Wuhan prior to the emergence of  
25 SARS-CoV-2, we collected bronchoalveolar lavage fluid samples from 408 patients presenting  
26 with pneumonia and acute respiratory infections at the Central Hospital of Wuhan between 2016  
27 and 2017. Unbiased total RNA sequencing was performed to reveal their “total infectome”,  
28 including viruses, bacteria and fungi. Consequently, we identified 37 pathogen species,  
29 comprising 15 RNA viruses, 3 DNA viruses, 16 bacteria and 3 fungi, often at high abundance  
30 and including multiple co-infections (12.8%). However, SARS-CoV-2 was not present. These  
31 data depict a stable core infectome comprising common respiratory pathogens such as  
32 rhinoviruses and influenza viruses, an atypical respiratory virus (EV-D68), and a single case of a  
33 sporadic zoonotic pathogen – *Chlamydia psittaci*. Samples from patients experiencing  
34 respiratory disease on average had higher pathogen abundance than healthy controls.  
35 Phylogenetic analyses of individual pathogens revealed multiple origins and global transmission  
36 histories, highlighting the connectedness of the Wuhan population. This study provides a  
37 comprehensive overview of the pathogens associated with acute respiratory infections and  
38 pneumonia, which were more diverse and complex than obtained using targeted PCR or qPCR  
39 approaches. These data also suggest that SARS-CoV-2 or closely related viruses were absent  
40 from Wuhan in 2016-2017.

41 **Key words:** Acute respiratory infection, pneumonia, total infectome, SARS-CoV-2, Wuhan

## 42 **Introduction**

43 The emergence of COVID-19 at the end of 2019 has had a profound impact on the world. The  
44 causative agent, a betacoronavirus (*Coronaviridae*) termed SARS-CoV-2, has high  
45 transmissibility, causes mild to severe respiratory symptoms in humans<sup>1-3</sup>. The first documented  
46 case of SARS-CoV-2 was reported in Wuhan, Hubei province, China<sup>3</sup>. Despite intensive research  
47 into SARS-CoV-2 and COVID-19, important questions remain unclear regarding the emergence  
48 of this virus, including for when and where it appeared and how long it was circulating in human  
49 populations prior to its initial detection in December 2019.

50 Meta-transcriptomics has several advantages over traditional diagnostic approaches based on  
51 serology or PCR<sup>4</sup>: it targets all types of micro-organisms simultaneously, identifies potential  
52 pathogens without *a priori* knowledge of what micro-organisms might be present, reveals the  
53 information (RNA) expressed by the pathogen during infection that is central to agent  
54 identification and studies of disease association. This method has been proven highly successful  
55 in revealing the entire virome and microbiome in a diverse range of species<sup>5-8</sup>, including the  
56 initial identification of SARS-CoV-2 from patients with severe pneumonia<sup>1</sup>.

57 Acute respiratory infections and pneumonia are a significant public health concern on a global  
58 scale. However, far less is known about the total “infectomes” associated with respiratory  
59 infections and pneumonia. Herein, we report total infectome surveillance of 408 patients  
60 presenting with pneumonia and acute respiratory infections at Wuhan Central Hospital prior to  
61 the SARS-CoV-2 epidemic. The purpose of this study was to use an un-biased meta-  
62 transcriptomics tool to characterize the total infectome within these patients. Nevertheless, since  
63 the sampling period occurred before the outbreak of COVID-19, this represents the first

64 opportunity to characterize the entire range of pathogens simultaneously within a cohort and  
65 determine the microbial composition of the population in which SARS-CoV-2 was initially  
66 reported.

67

## 68 **Results**

### 69 **Patient context**

70 We considered 408 patients clinically diagnosed with pneumonia or acute respiratory infection at  
71 the Central hospital of Wuhan in Wuhan, China. The sampling lasted for 20 months and covered  
72 the period between May 2016 to December 2017, two years before the onset of the COVID-19  
73 pandemic (Fig. 1A). The male-to-female ratio among the patients sampled was 1.4 (Fig. 1B),  
74 with age ranging from 16 to 90 years (medium, 62). Pre-existing medical conditions present in  
75 these patients included hypertension (n=108), diabetes (n=46), bronchiectasia (n=31), chronic  
76 obstructive pulmonary disease (COPD, n=23), cancer (n=12), and heart disease (n=10). Based on  
77 evaluations made by clinicians at the hospital, 27 patients were described as severely ill, with  
78 381 presenting with non-severe syndromes (Fig. 1E). The mortality for the entire cohort was  
79 0.74% (n=3) and the average duration of hospitalization was 8 days (range 2-322, medium 9).

### 80 **Total infectome**

81 Meta-transcriptomic analysis of the BALF samples identified a wide range of RNA viruses, DNA  
82 viruses, bacteria and fungi. For the purposes of this study, we only characterized those likely  
83 associated with human disease (i.e., pathogens). This included (i) existing species that are known  
84 to be associated with human disease, and (ii) potentially novel pathogens that have not been  
85 previously characterized. For the latter, we only considered DNA and RNA viruses that are  
86 related to a virus genus or family that have previously been shown to infect mammals and are at

87 relatively high abundance level (i.e., >0.1% of total RNA, or 1000 RPM). Other than new  
88 pathogens, the abundance threshold for pathogen positives was set as 1 RPM. Furthermore,  
89 commensal bacteria population was not considered here.

90 Based on these criteria we did not identify any potential novel viral pathogens. All the microbes  
91 identified belonged to those previously characterized as human pathogens, comprising 15 RNA  
92 viruses, 3 DNA viruses, 16 bacteria and 3 fungal pathogens (Fig. 2). The case positive rate for all  
93 pathogens was 71.25% (n=249, Fig. 2A), many of which were only associated with RNA viruses  
94 (27.3%, n=151) or bacteria (28.9%, n=160). Co-infection with two different pathogens was also  
95 commonplace, comprising a total of 71 (12.8%) cases (Fig. 2A). Among the pathogens  
96 identified, most were common respiratory pathogens such as influenza viruses, rhinoviruses,  
97 *Pseudomonas aeruginosa* and *Haemophilus influenzae* (Fig. 2B). In addition, we identified a  
98 number of unconventional respiratory pathogens that are often not included in respiratory  
99 pathogen screening panels but known to cause severe infections in respiratory tract or lungs,  
100 including enterovirus D68 and *Chlamydia psittaci* (see below).

101 Finally, none of the pathogens described here appeared in the blank controls. Since the blank  
102 control samples were generated using the same procedures for RNA extraction, library  
103 preparation and sequencing as the experimental groups, these results effectively exclude the  
104 possibility that the pathogens described above were of contaminant origin.

## 105 **Viruses**

106 RNA viral pathogens exhibited both great diversity (15 species) and abundance (up to 52% of  
107 total RNA) in the BALF samples examined here. The most frequently detected RNA viruses  
108 were human rhinoviruses A-C (HRV, n=55), followed by influenza A virus (IAV, n=29), human

109 parainfluenza virus type 3 (HPIV3, n=20), influenza B virus (IBV, n=8), and human  
110 metapneumovirus (HPMV, n=7) (Fig. 2B). While the majority were found throughout the study  
111 period, a few had more specific time-scales (Fig. 2C). For example, influenza B viruses were  
112 mostly identified in 2017, whereas enterovirus D (ENV-D, n=6) was only detected in the summer  
113 of 2016. In addition, we identified all four types of common cold associated coronaviruses -  
114 OC43 (n=4), HKU1 (n=4), 229E (n=6) and NL63 (n=1) - all of which had a relatively low  
115 prevalence in our cohort. Importantly, none of the libraries contained any hit to SARS-CoV or  
116 SARS-CoV-2, these results were confirmed by both read mapping and blasting against  
117 corresponding virus genomes.

118 In comparison to RNA viruses, the DNA viruses identified were limited in diversity and  
119 abundance. All three major types of human herpesviruses were identified - HSV1 (n=3), CMV  
120 (n=3), and EBV (n=2) - although with low abundance in all cases (up to 30 RPM or 0.003%)  
121 (Fig. 2C). Another common DNA virus that causes respiratory disease – adenovirus - was also  
122 identified in several cases, although at abundance levels lower than the 1 RPM threshold such  
123 that it was considered a ‘negative’ result in this context.

## 124 **Bacteria and fungi**

125 The most common bacterial pathogens identified included *Pseudomonas aeruginosa* (n=31),  
126 *Haemophilus influenzae* (n=19), *Staphylococcus aureus* (n=18), and *Mycoplasma pneumoniae*  
127 (n=11), all of which are common respiratory pathogens. *Acinetobacter* bacteria were also  
128 prevalent in our cohort, including *Acinetobacter baumannii* (n=5) and *A. pittii* (n=1), both of  
129 which are commonly associated with hospital-acquired infections. Other important respiratory  
130 pathogens, such as *Mycobacterium tuberculosis* (n=2), *Legionella pneumophila* (n=2),

131 *Streptococcus pneumoniae* (n=2), *Klebsiella pneumoniae* (n=3) and *Moraxella catarrhalis* (n=3),  
132 were also detected, although at a relatively low prevalence. Of particular interest was the  
133 identification of a single case of *Chlamydia psittaci* – a potentially bird-associated zoonotic  
134 pathogen – present in the BLAF at relatively high abundance (6396 RPM). Conversely, all the  
135 fungal pathogens identified here – *Candida albican* (n=3), *C. tropicalis* (n=1), and *Aspergillus*  
136 *spp.* (n=1) – were common pathogens known to cause respiratory infections.

### 137 **qPCR confirmation of pathogen presence and abundance**

138 For each of the pathogens identified here, we performed a qRT-PCR assay to confirm their  
139 presence and validate their abundance level as measured using our meta-transcriptomic  
140 approach. Strikingly, strong correlations were observed between the abundance measured by  
141 qPCR (i.e., CT value) and those estimated by read count after log 2 conversion ( $-0.8 < \text{Pearson's}$   
142  $R < -1$ , Fig. 3). Hence, the quantification by the two methods is strongly comparable.

### 143 **In-depth phylogenetic characterization of pathogens**

144 Although no novel viral pathogens were identified in this study, those viruses detected were  
145 characterized by substantial phylogenetic diversity, reflected in the presence of multiple viral  
146 lineages that highlight their complex epidemiological history in Wuhan (Fig. 4). We identified  
147 more than 14 genomic types of rhinovirus A, 7 of rhinovirus B, and 4 of rhinovirus C. A similar  
148 pattern of the co-circulation of multiple viral lineages was observed in other viruses. For  
149 example, the influenza A viruses identified in this study can be divided into the H1N1 and H3N2  
150 subtypes, each containing multiple lineages that clustered with viruses sampled globally and  
151 reflecting the highly connected nature of Wuhan (Fig. 4). In this context it was striking that in the  
152 case of several coronaviruses - OC43, HKU1 and 229E - the Wuhan sequences grouped directly

153 with those identified from the United States (Fig. 4), although this may reflect limited sampling.

#### 154 **Pathogen presence and abundance in diseased and healthy individuals**

155 Our meta-transcriptomic analysis revealed that many RNA viruses and bacteria detected were  
156 present at extremely high abundance levels (>1%, and up to 52% of total RNA) and hence likely  
157 indicative of acute disease. This was particularly true of eight species of RNA viruses – EV-D68,  
158 the influenza viruses, HRV, HPIV3, 229E – as well as two species of bacteria (*Haemophilus*  
159 *influenzae* and *Pseudomonas aeruginosa*) (Fig. 5). Together, these comprise a total of 54 cases  
160 (13.2% of total diseases cases).

161 In marked comparison, high levels of abundance were never observed in the healthy control  
162 group (Fig. 5). The highest abundance was 1302 RPM (0.013% of total RNA) in this group was  
163 for *Haemophilus influenzae*. Indeed, for most of the pathogens, the abundance level in healthy  
164 group was either undetectable or well below that observed in the diseased group, with the  
165 exception of Human parainfluenza virus type 1 (HPIV) and *Mycoplasma orale* for which the  
166 abundance levels were higher in the control group, although the sample size for both pathogens  
167 was relatively small. Conversely, commensal microbes, particularly *Escherichia coli*, exhibited  
168 similar abundance levels in the disease and control groups.

169

#### 170 **Discussion**

171 Our study provides a critical snapshot of the respiratory pathogens present in Wuhan prior to the  
172 emergence of SARS-CoV-2. Our unbiased metagenomic survey in patients presenting with  
173 pneumonia or acute respiratory infections provides strong evidence that SARS-CoV-2 or any  
174 related SARS-like viruses were absent in Wuhan approximately two years prior to the onset of



175 pandemic, although a number of common cold coronaviruses (HKU1, OC43, 229E, and NL63)  
176 were commonly detected in our cohort. Indeed, the earliest COVID-19 case, identified by qRT-  
177 PCR or next-generation sequencing-based assays performed at designated authoritative  
178 laboratories, can currently only be traced back to early December 2019 in Wuhan<sup>9</sup>. In addition, a  
179 retrospective survey of 640 throat swabs from patients with influenza-like illness in Wuhan from  
180 the period between 6 October 2019 and 21 January 2020 did not find any evidence of SARS-  
181 CoV-2 infection prior to January 2020<sup>10</sup>, such that the ultimate origin of SARS-CoV-2 remains  
182 elusive<sup>11</sup>.

183 The data presented provide a comprehensive overview of the infectome associated with  
184 pneumonia or acute respiratory infections in Wuhan, which is clearly more diverse and complex  
185 than described using previous surveys based on targeted PCR or qPCR approach alone<sup>12</sup>. In light  
186 of our observations, we can divide the respiratory infectome into three categories based on  
187 epidemiological characteristics: (i) the “core” infectome that is commonly found in patients with  
188 respiratory infection and is expected to occur each year globally, (ii) an “emerging” infectome  
189 that occurs during outbreaks but are not typically found in the geographic regions under  
190 investigation, and (iii) the sporadic infection or zoonotic infectome of new or rare pathogens.

191 The core infectome comprised of a wide range of common respiratory or systemic pathogens that  
192 are subject to frequent screening in hospitals. These include influenza viruses, HMPV, RSV,  
193 *Moraxella catarrhalis*, *Acinetobacter spp.*, *Klebsiella pneumoniae*, *Mycoplasma spp.*,  
194 *Haemophilus influenzae*, *Pseudomonas aeruginosa*, *Staphylococcus aureus* and *Streptococcus*  
195 *pneumoniae*. However, the remaining pathogens identified here, including rhinoviruses,  
196 parainfluenza viruses, coronaviruses, and herpesviruses, have often received far less attention

197 from clinicians and are sometimes ignored entirely, most likely due to the lack of association  
198 with severe disease in adults<sup>13-15</sup>. Nevertheless, our results showed these “neglected” respiratory  
199 viruses had high diversity, abundance and prevalence in the cohort of pneumonia or acute  
200 respiratory patients studied here in comparison to healthy controls, such that their role as agents  
201 of disease should not be underestimated. One scenario is that they represent opportunistic  
202 pathogens that take advantage of weakened immunity, such as herpesviruses associated acute  
203 respiratory distress (ARDS)<sup>16</sup>. It is also possible that their pathogenic effects have yet to be  
204 identified and may extend to disease manifestations beyond respiratory infections. For example,  
205 deep sequencing of a brain biopsy sample suggesting that the OC43 coronavirus may be  
206 associated with fatal encephalitis in humans<sup>17</sup>.

207 We also identified a potential “emerging” infectome, in this case comprising a single virus - EV-  
208 D68 - that may represent a regional or national outbreak of an unconventional respiratory  
209 pathogens. The prevalence of EV-D68 remained low from its discovery in 1967 until 2014<sup>18</sup>,  
210 when a major outbreak started in the United States and spread to more than 20 countries<sup>19</sup>,  
211 causing severe respiratory illness with potential neurological manifestations such as acute flaccid  
212 paralysis in children<sup>20, 21</sup>. In China, EV-D68 was only sporadically reported<sup>22, 23</sup>, although  
213 serological surveys suggest a much wider prevalence for both children and adults since 2009<sup>24</sup>.  
214 We identified six EV-D68 cases, all adults that occurred within a relatively narrow time window  
215 between June and December 2016. Phylogenetic analysis revealed that the sequences of these  
216 viruses were closely related to each other and to other Chinese strains from the same period (Fig.  
217 4), suggesting that it may be a part of a larger outbreak in China. The EV-D68 cases identified  
218 here showed moderate to severe respiratory symptoms, although viral abundance was generally  
219 very high, with four of six cases showing >106 RPM (i.e., >10% of total RNA) in the BALF

220 sample. This highlights the active replication and massive proliferation of viruses within the  
221 respiratory system of these patients.

222 Finally, our zoonotic infectome also comprised a single pathogen, *Chlamydia psittaci*, that is  
223 associated with avian species but causes occasional outbreaks in domestic animals (i.e., pigs,  
224 cattles, and sheep) and humans<sup>25</sup>. In humans, *C. psittaci* infections often starts with influenza-  
225 like symptoms, but can develop into serious lung infections and even death<sup>26</sup>. The single case of  
226 *C. psittaci* identified here was at relatively high abundance level (6396 RPM) and caused a  
227 relatively severe disease, with the patient experiencing expiratory dyspnea, severe pneumonia  
228 and pleural effusion, and was subsequently transfer to intensive care unit (ICU) for further  
229 treatment. Since the patient had no travel history for a month prior to illness, this discovery  
230 underlines the risk of local exposure to this bacterium.

231 Our phylogenetic analysis of the metagenomic data generated revealed extensive intra-specific  
232 diversity in each virus species identified highlights the complex epidemiological history of these  
233 pathogens. Indeed, the influenza A viruses (H1N1 and H3N2), influenza B virus, HPIV3 and  
234 HCoV-OC43 discovered here all comprised multiple lineages (Fig. 4), suggesting these viruses  
235 were introduced from diverse sources. Since some of the viruses were closely related to those  
236 circulating in other countries it is possible that they represent overseas importations: this is not  
237 surprising given that Wuhan is a major domestic travel hub and well linked internationally.  
238 Indeed, the rapid and widespread transmission of SARS-CoV-2 between Wuhan and other major  
239 cities globally was key to seeding the global pandemic of COVID-19, with early cases in a  
240 number of localities all linked to travel from Wuhan<sup>26-29</sup>.

241 The methodology used here served as an unbiased investigation of potential emerging pathogens.

242 Meta-transcriptomics has several advantages over traditional diagnostic approaches based on  
243 serology or PCR: it targets all types of micro-organisms simultaneously, identifies potential  
244 pathogens without *a priori* knowledge of what micro-organisms might be present, reveals the  
245 information (RNA) expressed by the pathogen during infection that is central to agent  
246 identification and studies of disease association. With the popularization of next-generation  
247 sequencing platform in major hospitals, the approach outlined here can be easily integrated into  
248 diagnostic practice with much greater speed and significantly more information output than  
249 traditional technologies, providing a broad-scale understanding of infectious disease in general.

250

## 251 **Material and Methods**

### 252 **Ethics statement**

253 The sampling and experimental procedures for this study were reviewed by the ethics  
254 committees of the Central Hospital of Wuhan and the National Institute for Communicable  
255 Disease Control and Prevention, Chinese Center for Disease Control and Prevention. Written  
256 informed consents were taken from all patients and volunteers recruited in this study. In addition,  
257 for child patients, written consents were also obtained from their parents or guardians. Physicians  
258 were informed of results of the pathogen discovery exercise as soon as the meta-transcriptomic  
259 results were obtained.

260

### 261 **Sample collection from patients and controls.**

262 More than 1000 patients were recruited from the Central Hospital of Wuhan between 2016 and  
263 2017. The target clinical conditions were community-acquired pneumonia and acute respiratory  
264 infections based on the initial diagnosis made by clinicians. All patients were hospitalized and

265 subject to bronchoalveolar lavage fluid (BALF) collection required by the initial diagnosis for  
266 pneumonia or acute respiratory distress syndrome and independent of this study. The BALF  
267 sample was divided into two parts for the clinical laboratory test and this study, respectively. Of  
268 the BALF samples collected, 409 were subjected to meta-transcriptomic analysis based on their  
269 condition and the time between hospitalization and sample collection. No diagnostic information  
270 was provided prior to sample selection. To establish a healthy control group, 5ml-10ml of BALF  
271 samples were also taken from 32 volunteers without respiratory symptoms. We also included 10  
272 blank controls where only RNase free water was used for nucleic acid extraction and library  
273 construction, although only four of these produced viable RNA sequencing results.

#### 274 **Meta-transcriptomic pathogen discovery pipeline**

275 We followed a standard protocol for meta-transcriptomics analysis for each BALF sample. Total  
276 RNA was first extracted from 200–300ul of each sample using the RNeasy Plus Universal Kit  
277 (Qiagen, USA) according to the manufacturer’s instructions. From the extracted RNA, we  
278 performed human rRNA removal and low concentration library construction procedures with the  
279 Trio RNA-Seq kit (NuGEN Technologies, USA). The libraries were then subjected to 150bp  
280 pair-end sequencing on an Illumina Hiseq 4000 platform at Novogene (Beijing), with target  
281 output of 10G base pairs per library. For each of the sequencing results generated, we removed  
282 adaptor sequences, non-complex reads, as well as duplicated reads using the BBmap software  
283 package. Human and ribosomal RNA (rRNA) reads were subsequently removed by mapping the  
284 de-duplicated reads against the human reference genome (GRCh38/hg38) and the comprehensive  
285 rRNA sequence collection downloaded from SILVA database<sup>30</sup>.

286 The remaining sequencing reads were subject to a pathogen discovery pipeline. For virus

287 identification, the reads were directly compared against reference virus databases using the  
288 blastn program and against the non-redundant protein (nr) database using diamond blastx<sup>31</sup>, with  
289 an e-value threshold set at 1E-10 and 1E-5 for blastn and diamond blastx analyses, respectively.  
290 Viral abundance was summarized from both analyses, calculated using the relation: total viral  
291 reads/total non-redundant reads\* 1 million (i.e., reads per million, RPM). To identify highly  
292 divergent virus genomes, reads were assembled using megahit<sup>32</sup> into contigs before comparison  
293 against the nt and nr databases. Those reads with <90% amino acid similarity to known viruses  
294 were treated as potential novel virus species. For bacterial and fungi identification, we first used  
295 MetaPhlan2<sup>33</sup> to identify potential species in both groups. Relevant background bacterial and  
296 fungal genomes were subsequently downloaded from NCBI/GenBank and used as a template for  
297 read mapping. Based on the mapping results for each case, we generated relevant contigs for  
298 blastn analyses against the non-redundant nucleotide (nt) database to determine taxonomy to the  
299 species level. The abundance level of bacterial and fungi pathogens was also calculated in the  
300 form of RPM based on genome and mitochondrial genome read counts, respectively.

301 A microbe was considered as “positive” within a specific sample if its abundance level was  
302 greater than 1 RPM. To prevent false positives resulting from index hopping, we used a threshold  
303 of 0.1% for viruses present in the same sequencing lane: that is, if the libraries contain less than  
304 0.1% of the most abundant library it is treated as “negative”.

### 305 **Confirmatory testing by conventional methods**

306 For pathogen positive samples, the same RNA used for meta-transcriptomics analysis was also  
307 subject to a qRT-PCR assay with primers sets designed for a specific or related group of  
308 pathogens (Table S1). RNA was first reverse transcribed by SuperScript™ III First-Strand

309 Synthesis SuperMix for qRT-PCR (Invitrogen, California), and then amplified by TaqPath™  
310 ProAmp™ Master Mix (Applied Biosystems, California). A cycle threshold (CT) value of 38 and  
311 above was treated as negative.

### 312 **Pathogen genomic analyses**

313 For viruses at high abundance levels (i.e., > 1000 RPM), complete genomes were assembled  
314 using Megahit and confirmed by mapping reads against the assembled contigs. These genomes  
315 were then aligned with related reference virus sequences downloaded from NCBI/GenBank  
316 using MAFFT program<sup>34</sup>. Ambiguously aligned regions were removed using the Trimal  
317 program<sup>35</sup>. Phylogenetic trees were estimated using the maximum likelihood approach  
318 implemented in PhyML<sup>36</sup>, employing the GTR model of nucleotide substitution and SPR branch  
319 swapping. The support for each node in the tree was estimated with an approximate likelihood  
320 ratio test (aLRT) with the Shimodaira-Hasegawa-like procedures.

### 321 **Data availability**

322 All non-human reads have been deposited in the SRA databases under the project accession  
323 PRJNA699976. Relevant consensus virus genome sequences have been deposited in  
324 NCBI/GenBank under the accessions MW567157- MW567162, MW570805- MW570808,  
325 MW571087- MW571107, MW587035- MW587095.

326

327

### 328 **Acknowledgments**

329 This study was supported by the National Natural Science Foundation of China 32041004  
330 (YZZ), 31930001 (YZZ), 81861138003 (YZZ) and 81672057 (YZZ). E.C.H. is supported by an

331 ARC Australian Laureate Fellowship (FL170100022).

332

333 **Author Contributions**

334 Yong-Zhen Zhang conceived and designed the study. Su Zhao, Yi Hu, Wen Yin, Fang Ni, Hong-

335 Ling Hu, Shuang Geng, and Li Tan performed the clinical work and sample collection. Bin Yu,

336 Jun-Hua Tian and Ying Peng performed epidemiological investigation and sample collection.

337 Mang Shi, Wei-Chen Wu, Yan-Mei Chen, Wen Wang, and Zhi-Gang Song performed the

338 experimental works. Mand Shi, Edwards C Holmes and Yong-Zhen Zhang analysed the data.

339 Mang Shi, Edward C Holmes and Yong-Zhen Zhang wrote the paper with input from all authors.

340

341 **Competing Interests**

342 All other authors declare no competing interests.



## 343 References

- 344 1. Wu F, Zhao S, Yu B, Chen YM, Wang W, Song ZG, Hu Y, Tao ZW, Tian JH, Pei YY, Yuan  
345 ML, Zhang YL, Dai FH, Liu Y, Wang QM, Zheng JJ, Xu L, Holmes EC, Zhang YZ. A new  
346 coronavirus associated with human respiratory disease in China. *Nature* 2020; 579:265-9.
- 347 2. Zhou P, Yang XL, Wang XG, Hu B, Zhang L, Zhang W, Si HR, Zhu Y, Li B, Huang CL,  
348 Chen HD, Chen J, Luo Y, Guo H, Jiang RD, Liu MQ, Chen Y, Shen XR, Wang X, Zheng  
349 XS, Zhao K, Chen QJ, Deng F, Liu LL, Yan B, Zhan FX, Wang YY, Xiao GF, Shi ZL. A  
350 pneumonia outbreak associated with a new coronavirus of probable bat origin. *Nature* 2020;  
351 579:270-3.
- 352 3. Li Q, Guan X, Wu P, Wang X, Zhou L, Tong Y, Ren R, Leung KSM, Lau EHY, Wong JY,  
353 Xing X, Xiang N, Wu Y, Li C, Chen Q, Li D, Liu T, Zhao J, Liu M, Tu W, Chen C, Jin L,  
354 Yang R, Wang Q, Zhou S, Wang R, Liu H, Luo Y, Liu Y, Shao G, Li H, Tao Z, Yang Y, Deng  
355 Z, Liu B, Ma Z, Zhang Y, Shi G, Lam TTY, Wu JT, Gao GF, Cowling BJ, Yang B, Leung  
356 GM, Feng Z. Early transmission dynamics in Wuhan, China, of novel coronavirus-infected  
357 pneumonia. *N Engl J Med* 2020; 382:1199-207.
- 358 4. Zhang YZ, Chen YM, Wang W, Qin XC, Holmes EC. Expanding the RNA virosphere by  
359 unbiased metagenomics. *Annu Rev Virol*; 2019; 6:119-139.
- 360 5. Eden JS, Rose K, Ng J, Shi M, Wang Q, Sintchenko V, Holmes EC. *Francisella tularensis*  
361 *ssp. holarctica* in Ringtail Possums, Australia. *Emerg Infect Dis* 2017; 23:1198-201.
- 362 6. Shi M, Lin XD, Chen X, Tian JH, Chen LJ, Li K, Wang W, Eden JS, Shen JJ, Liu L, Holmes  
363 EC, Zhang YZ. The evolutionary history of vertebrate RNA viruses. *Nature* 2018; 556:197-  
364 202.
- 365 7. Shi M, Lin XD, Tian JH, Chen LJ, Chen X, Li CX, Qin XC, Li J, Cao JP, Eden JS,  
366 Buchmann J, Wang W, Xu JG, Holmes EC, Zhang YZ. Redefining the invertebrate RNA  
367 virosphere. *Nature* 2016; 540:539-43.
- 368 8. Li CX, Li W, Zhou J, Zhang B, Feng Y, Xu CP, Lu YY, Holmes EC, Shi M. High resolution  
369 metagenomic characterization of complex infectomes in paediatric acute respiratory

- 370 infection. *Sci Rep* 2020; 10:3963.
- 371 9. Huang C, Wang Y, Li X, Ren L, Zhao J, Hu Y, Zhang L, Fan G, Xu J, Gu X, Cheng Z, Yu T,  
372 Xia J, Wei Y, Wu W, Xie X, Yin W, Li H, Liu M, Xiao Y, Gao H, Guo L, Xie J, Wang G,  
373 Jiang R, Gao Z, Jin Q, Wang J, Cao B. Clinical features of patients infected with 2019 novel  
374 coronavirus in Wuhan, China. *Lancet* 2020; 395:497-506.
- 375 10. Kong WH, Li Y, Peng MW, Kong DG, Yang XB, Wang L, Liu MQ. SARS-CoV-2 detection  
376 in patients with influenza-like illness. *Nat Microbiol* 2020; 5:675-8.
- 377 11. Hu B, Guo H, Zhou P, Shi ZL. Characteristics of SARS-CoV-2 and COVID-19. *Nat Rev*  
378 *Microbiol* 2020.
- 379 12. Wu Z, Li Y, Gu J, Zheng H, Tong Y, Wu Q. Detection of viruses and atypical bacteria  
380 associated with acute respiratory infection of children in Hubei, China. *Respirology* 2014;  
381 19:218-24.
- 382 13. Jacobs SE, Lamson DM, St George K, Walsh TJ. Human rhinoviruses. *Clin Microbiol Rev*  
383 2013;26:135-62.
- 384 14. Henrickson KJ. Parainfluenza viruses. *Clin Microbiol Rev* 2003;16:242-64.
- 385 15. van der Hoek L. Human coronaviruses: what do they cause? *Antivir Ther* 2007;12:651-8.
- 386 16. Bonizzoli M, Arvia R, di Valvasone S, Liotta F, Zakrzewska K, Azzi A, Peris A. Human  
387 herpesviruses respiratory infections in patients with acute respiratory distress (ARDS). *Med*  
388 *Microbiol Immunol* 2016; 205:371-9.
- 389 17. Morfopoulou S, Brown JR, Davies EG, Anderson G, Virasami A, Qasim W, Chong WK,  
390 Hubank M, Plagnol V, Desforages M, Jacques TS, Talbot PJ, Breuer J. Human coronavirus  
391 OC43 associated with fatal encephalitis. *N Engl J Med* 2016; 375:497-8.
- 392 18. Schieble JH, Fox VL, Lennette EH. A probable new human picornavirus associated with  
393 respiratory diseases. *Am J Epidemiol* 1967; 85:297-310.
- 394 19. Holm-Hansen CC, Midgley SE, Fischer TK. Global emergence of enterovirus D68: a  
395 systematic review. *Lancet Infect Dis* 2016; 16:e64-e75.
- 396 20. Midgley CM, Watson JT, Nix WA, Curns AT, Rogers SL, Brown BA, Conover C,

- 397 Dominguez SR, Feikin DR, Gray S, Hassan F, Hoferka S, Jackson MA, Johnson D, Leshem  
398 E, Miller L, Nichols JB, Nyquist AC, Obringer E, Patel A, Patel M, Rha B, Schneider E,  
399 Schuster JE, Selvarangan R, Seward JF, Turabelidze G, Oberste MS, Pallansch MA, Gerber  
400 SI; EV-D68 Working Group. Severe respiratory illness associated with a nationwide  
401 outbreak of enterovirus D68 in the USA (2014): a descriptive epidemiological investigation.  
402 *Lancet Respir Med* 2015; 3:879-87.
- 403 21. Greninger AL, Naccache SN, Messacar K, Clayton A, Yu G, Somasekar S, Federman S,  
404 Stryke D, Anderson C, Yagi S, Messenger S, Wadford D, Xia D, Watt JP, Van Haren K,  
405 Dominguez SR, Glaser C, Aldrovandi G, Chiu CY. A novel outbreak enterovirus D68 strain  
406 associated with acute flaccid myelitis cases in the USA (2012-14): a retrospective cohort  
407 study. *Lancet Infect Dis* 2015;15:671-82.
- 408 22. Xiang Z, Gonzalez R, Wang Z, Ren L, Xiao Y, Li J, Li Y, Vernet G, Paranhos-Baccalà G, Jin  
409 Q, Wang J. Coxsackievirus A21, enterovirus 68, and acute respiratory tract infection, China.  
410 *Emerg Infect Dis* 2012;18:821-4.
- 411 23. Xiao Q, Ren L, Zheng S, Wang L, Xie X, Deng Y, Zhao Y, Zhao X, Luo Z, Fu Z, Huang A,  
412 Liu E. Prevalence and molecular characterizations of enterovirus D68 among children with  
413 acute respiratory infection in China between 2012 and 2014. *Sci Rep* 2015; 5:16639.
- 414 24. Xiang Z, Li L, Ren L, Guo L, Xie Z, Liu C, Li T, Luo M, Paranhos-Baccalà G, Xu W, Wang J.  
415 Seroepidemiology of enterovirus D68 infection in China. *Emerg Microbes Infect* 2017;  
416 6:e32.
- 417 25. Knittler MR, Sachse K. Chlamydia psittaci: update on an underestimated zoonotic agent.  
418 *Pathog Dis* 2015; 73:1-15.
- 419 26. Watanabe M. The COVID-19 Pandemic in Japan. *Surg Today* 2020; 50:787-93.
- 420 27. Lescure FX, Bouadma L, Nguyen D, Parisey M, Wicky PH, Behillil S, Gaymard A,  
421 Bouscambert-Duchamp M, Donati F, Le Hingrat Q, Enouf V, Houhou-Fidouh N, Valette M,  
422 Mailles A, Lucet JC, Mentre F, Duval X, Descamps D, Malvy D, Timsit JF, Lina B, van-der-  
423 Werf S, Yazdanpanah Y.. Clinical and virological data of the first cases of COVID-19 in  
424 Europe: a case series. *Lancet Infect Dis* 2020; 20:697-706.
- 425 28. Wei WE, Li Z, Chiew CJ, Yong SE, Toh MP, Lee VJ. Presymptomatic transmission of

- 426 SARS-CoV-2 - Singapore, January 23-March 16, 2020. *MMWR Morb Mortal Wkly Rep*  
427 2020; 69:411-5.
- 428 29. Song JY, Yun JG, Noh JY, Cheong HJ, Kim WJ. Covid-19 in South Korea - Challenges of  
429 subclinical manifestations. *N Engl J Med* 2020; 382:1858-9.
- 430 30. Quast C, Pruesse E, Yilmaz P, Gerken J, Schweer T, Yarza P, Peplies J, Glöckner FO. The  
431 SILVA ribosomal RNA gene database project: improved data processing and web-based  
432 tools. *Nucleic Acids Res* 2013; 41:D590-6.
- 433 31. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND. *Nat*  
434 *Methods* 2015; 12:59-60.
- 435 32. Li D, Liu CM, Luo R, Sadakane K, Lam TW. MEGAHIT: an ultra-fast single-node solution  
436 for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*  
437 2015; 31:1674-6.
- 438 33. Truong DT, Franzosa EA, Tickle TL, Scholz M, Weingart G, Pasolli E, Tett A, Huttenhower  
439 C, Segata N.. MetaPhlan2 for enhanced metagenomic taxonomic profiling. *Nat Methods*  
440 2015; 12:902-3.
- 441 34. Katoh K, Toh H. Recent developments in the MAFFT multiple sequence alignment program.  
442 *Brief Bioinform* 2008; 9:286-98.
- 443 35. Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. trimAl: a tool for automated alignment  
444 trimming in large-scale phylogenetic analyses. *Bioinformatics* 2009; 25:1972-3.
- 445 36. Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms  
446 and methods to estimate maximum-likelihood phylogenies: assessing the performance of  
447 PhyML 3.0. *Syst Biol* 2010; 59:307-21.  
448

449 **Figure legends**

450 **Fig. 1. Patient recruitment.** (A) Sampling frequency and intensity during the study period. (B)  
451 Male-to-female ratio of the recruited patients. (C) Age structure of all patients involved in this  
452 study depicted using violin plots. (D) Type and frequencies of pre-existing conditions. (E) The  
453 number of severe, non-severe and control cases.

454

455 **Fig. 2. Prevalence and abundance of viral, bacterial and fungal pathogens in the 408 cases**  
456 **examined in this study.** (A) Proportion of cases infected with RNA virus, DNA virus, bacteria,  
457 fungus and with mixed infections. (B) Prevalence of each pathogen, ordered by the number of  
458 cases. (C) Heat map showing the prevalence and abundance of pathogens in diseased and control  
459 samples. *Escherichia coli* is not regarded as pathogen but is shown here as an example of non-  
460 pathogen contamination. The samples (x-axis) are divided into “Patient” and “Control” groups,  
461 each ordered chronically. The pathogens (y-axis) are divided into four categories: RNA viruses,  
462 DNA viruses, bacteria, and fungi.

463

464 **Fig. 3. Comparisons of pathogen abundance measured by qPCR and meta-transcriptomics**  
465 **approaches.** Abundance by qPCR methods is measured by cycle threshold, or CT value, while  
466 those by meta-transcriptomics are measured by “log<sub>2</sub>(number of reads)”. The comparisons are  
467 performed on four most abundant pathogens: influenza A virus, influenza B virus, human  
468 rhinoviruses, and human parainfluenza virus 3. Pearson’s correlation coefficient between CT  
469 values and log<sub>2</sub>(number of reads) is estimated for each pathogen.

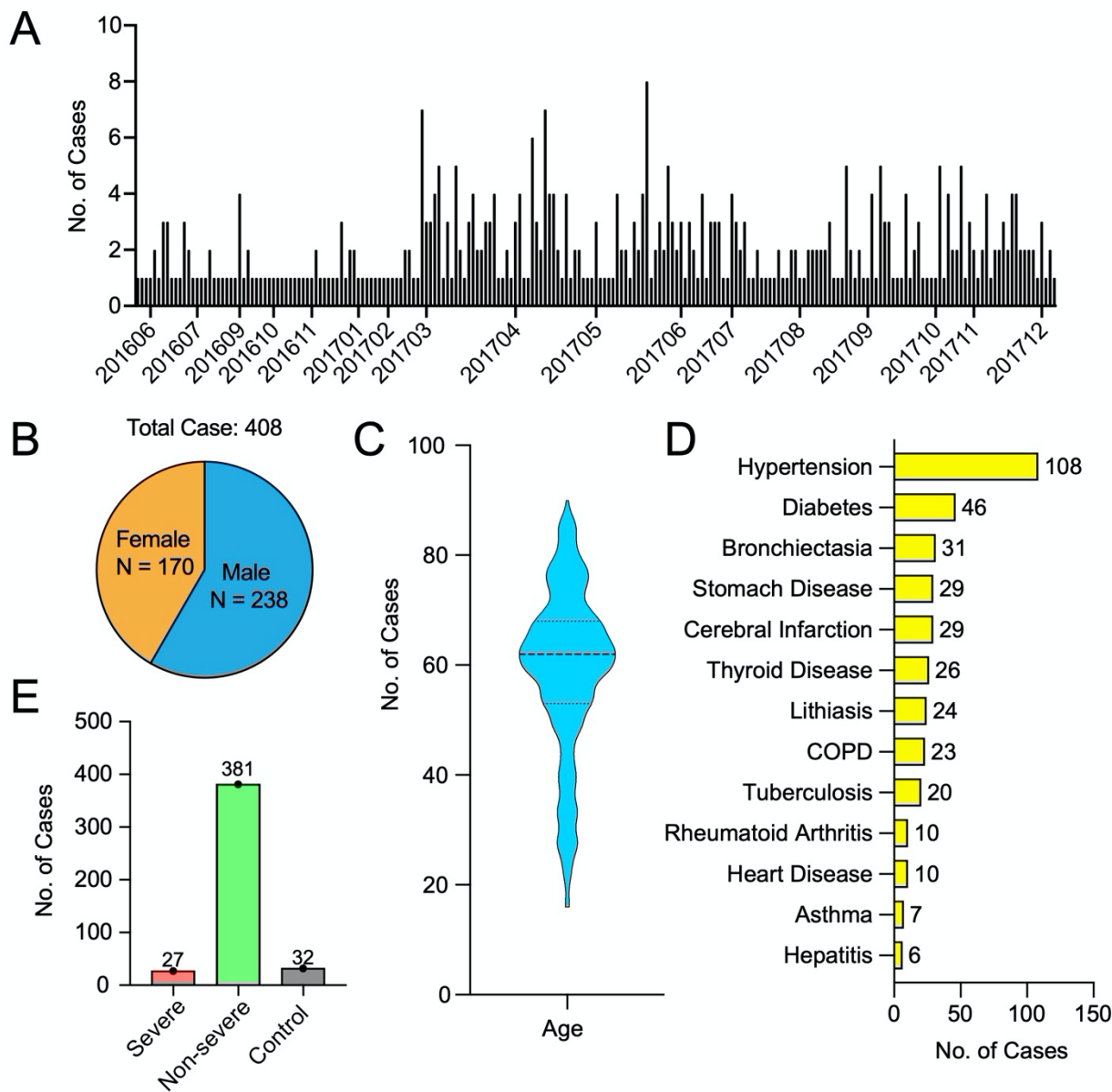
470

471 **Fig. 4. Evolutionary relationships of RNA virus pathogens identified in this study.**

472 Phylogenetic trees were estimated using the maximum likelihood method implemented in  
473 PhyML. Sequences identified from this study were marked with red solid circle. For larger trees,  
474 we only show lineages or sub-lineages which contain sequences identified in this study.

475

476 **Fig. 5. Comparison of prevalence levels between healthy and control groups.** Boxplots for  
477 patient and control groups for each pathogen identified in this study, including RNA virus  
478 (green), DNA virus (yellow), bacteria (orange) and fungus (blue). For clarity, only non-zero  
479 abundance levels are reported.



480

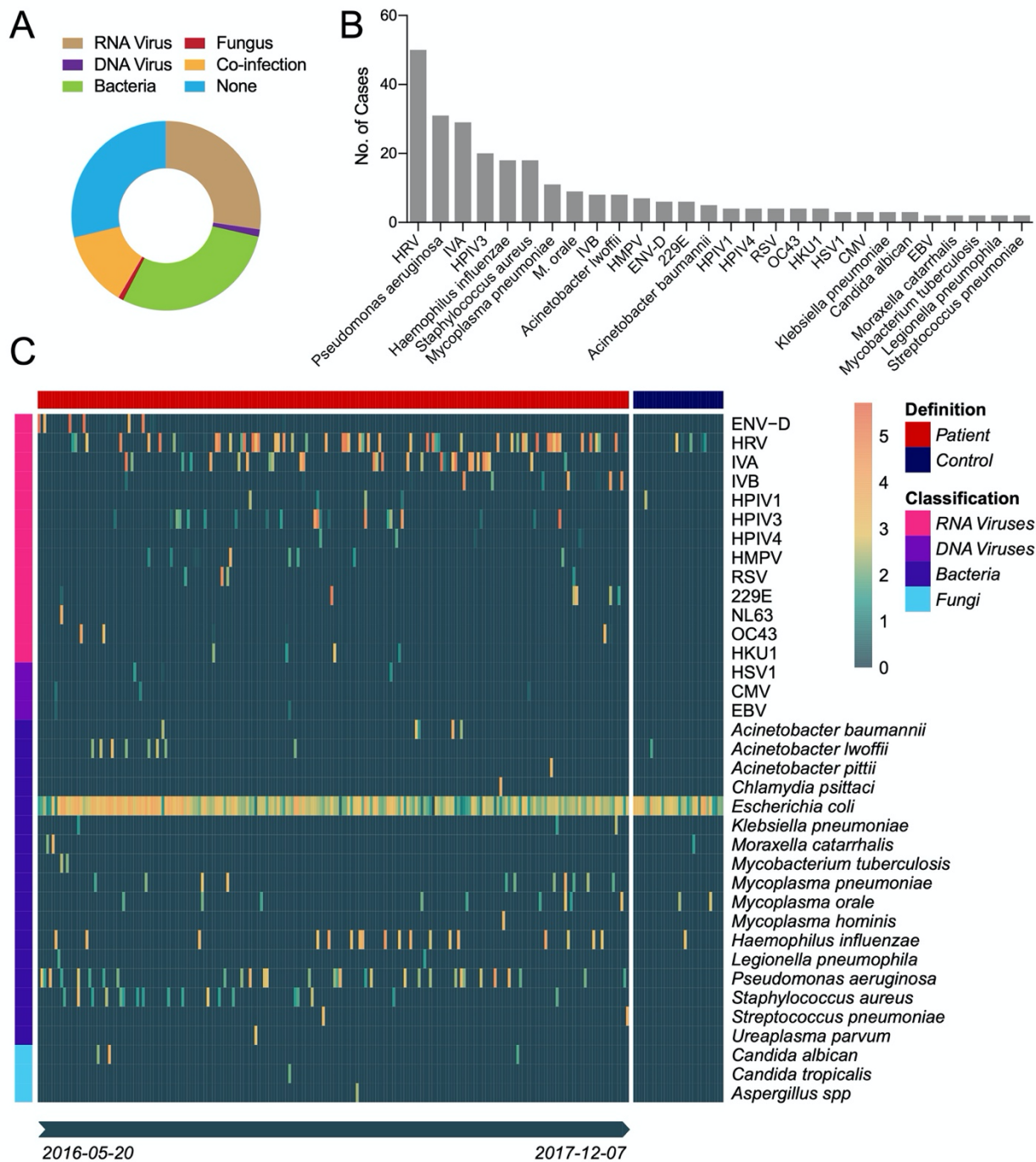
481 **Fig. 1. Patient recruitment.** (A) Sampling frequency and intensity during the study period. (B)

482 Male-to-female ratio of the recruited patients. (C) Age structure of all patients involved in this

483 study depicted using violin plots. (D) Type and frequencies of pre-existing conditions. (E) The

484 number of severe, non-severe and control cases.

485



486

487 **Fig. 2. Prevalence and abundance of viral, bacterial and fungal pathogens in the 408 cases**

488 **examined in this study.** (A) Proportion of cases infected with RNA virus, DNA virus, bacteria,

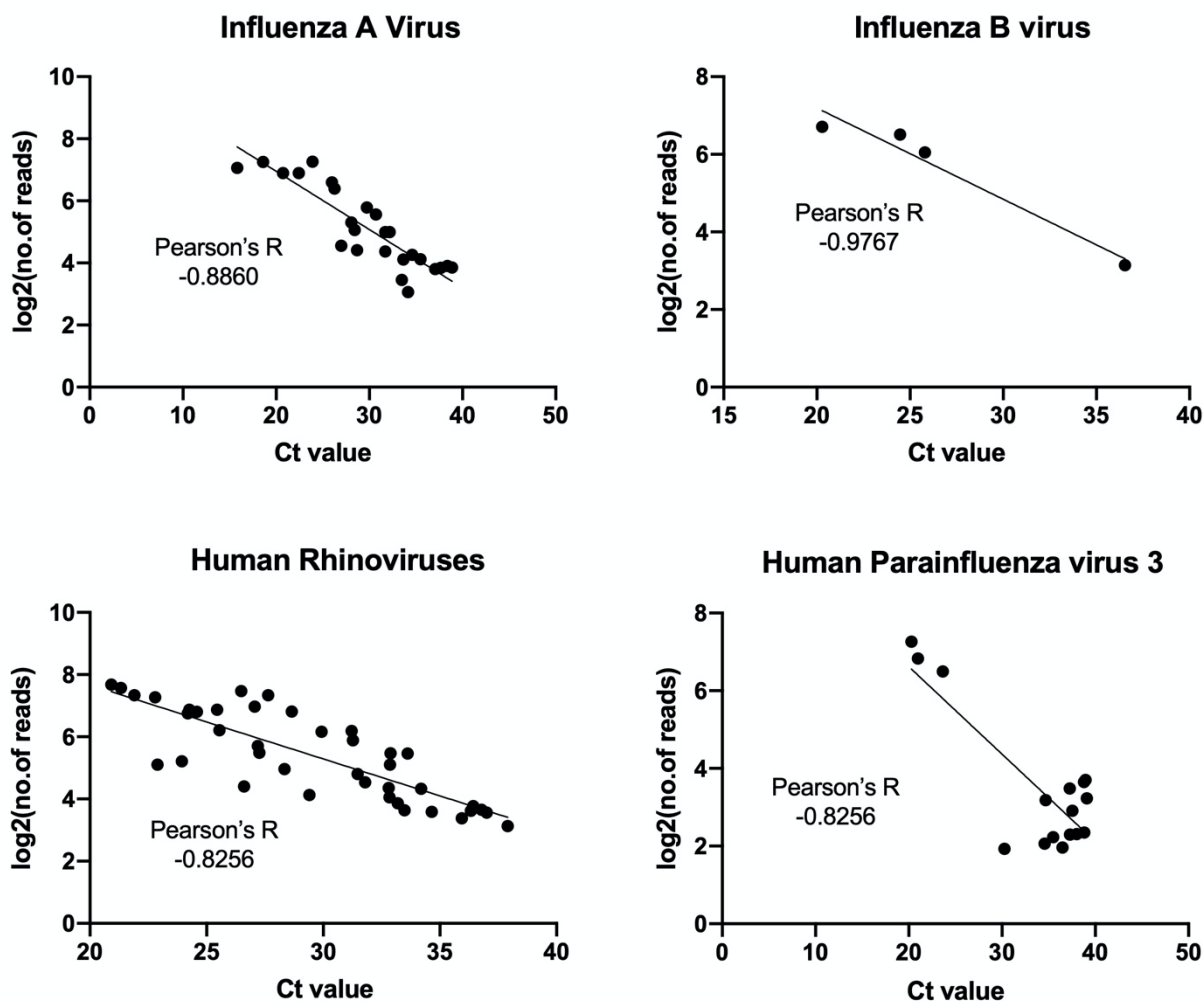
489 fungus and with mixed infections. (B) Prevalence of each pathogen, ordered by the number of

490 cases. (C) Heat map showing the prevalence and abundance of pathogens in diseased and control



491 samples. *Escherichia coli* is not regarded as pathogen but is shown here as an example of non-  
492 pathogen contamination. The samples (x-axis) are divided into “Patient” and “Control” groups,  
493 each ordered chronically. The pathogens (y-axis) are divided into four categories: RNA viruses,  
494 DNA viruses, bacteria, and fungi.

495



496

497 **Fig. 3. Comparisons of pathogen abundance measured by qPCR and meta-transcriptomics**

498 **approaches.** Abundance by qPCR methods is measured by cycle threshold, or CT value, while

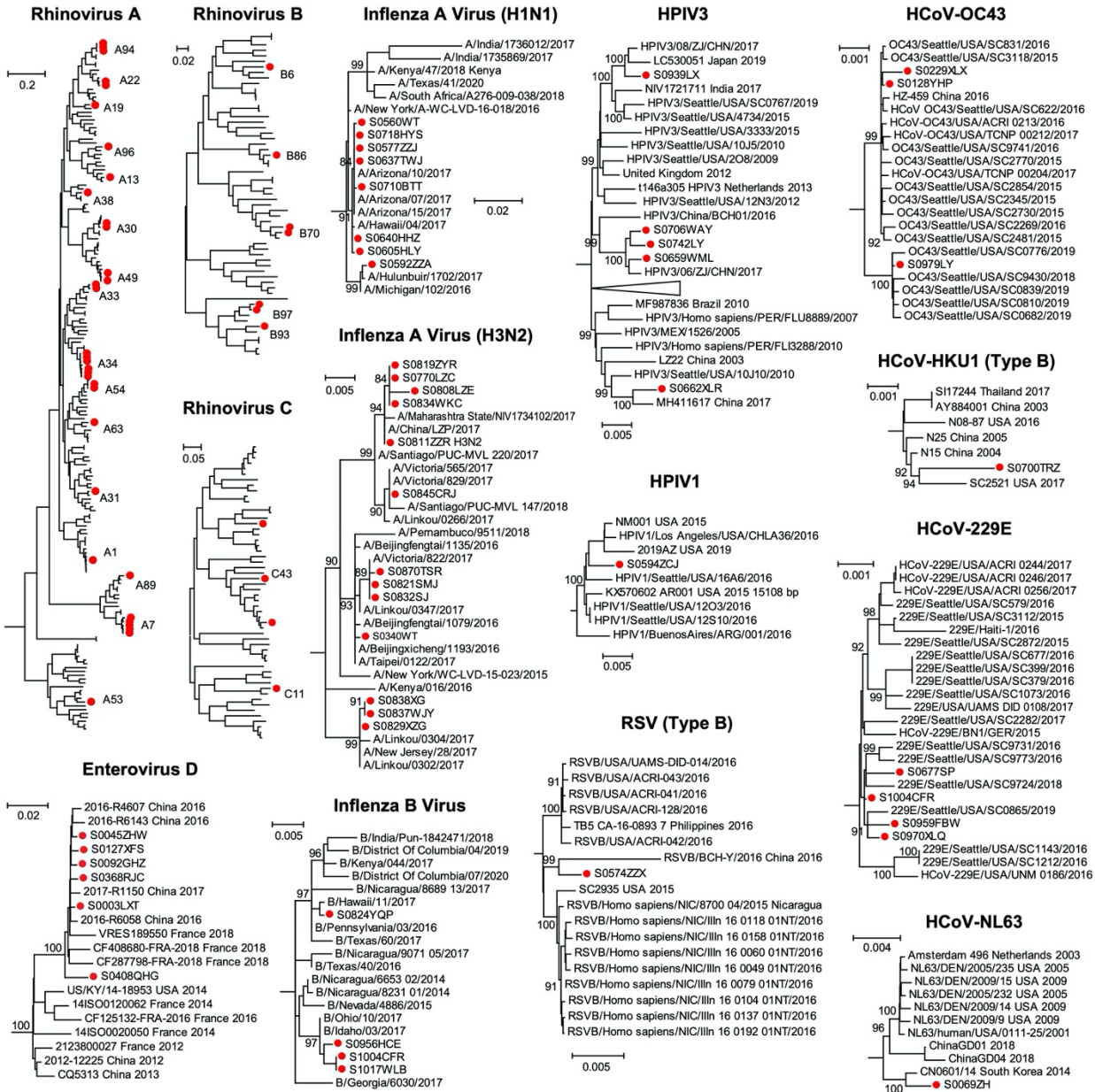
499 those by meta-transcriptomics are measured by “log2(number of reads)”. The comparisons are

500 performed on four most abundant pathogens: influenza A virus, influenza B virus, human

501 rhinoviruses, and human parainfluenza virus 3. Pearson’s correlation coefficient between CT

502 values and log2(number of reads) is estimated for each pathogen.

503



504

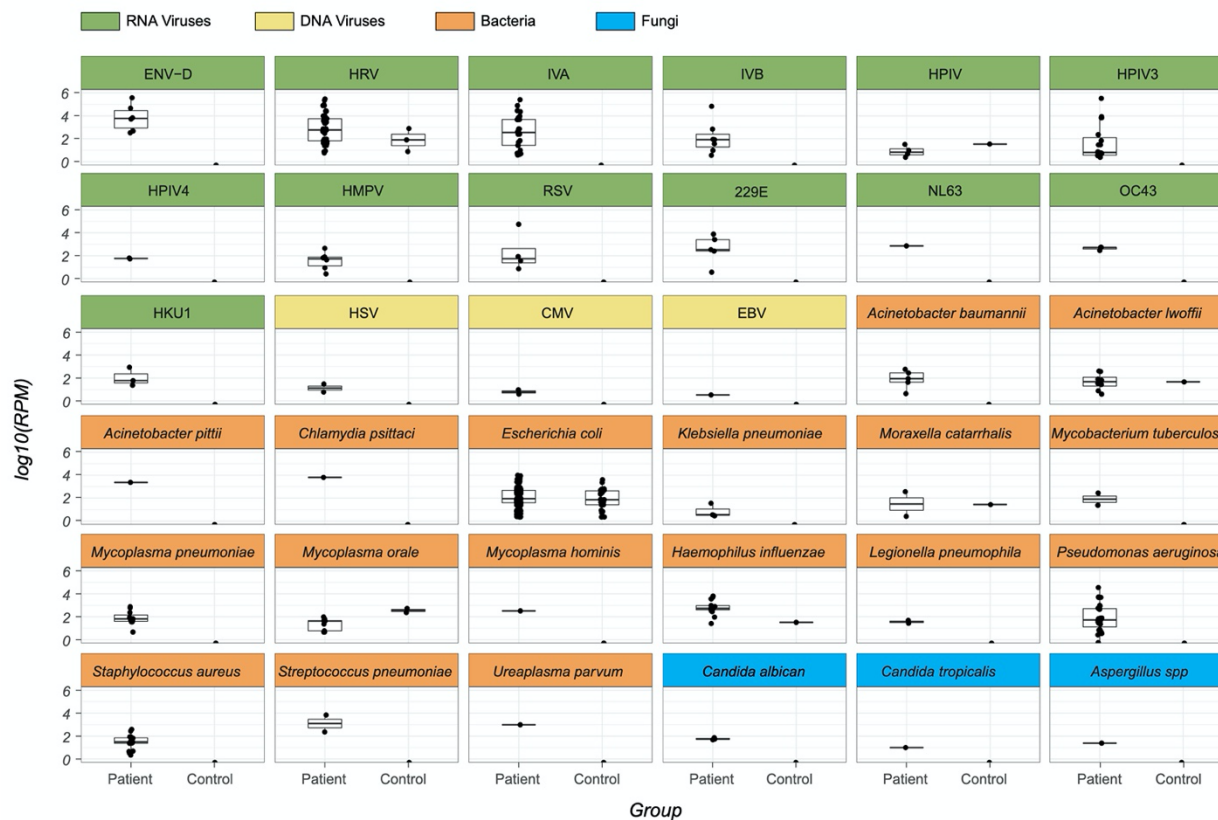
505 **Fig. 4. Evolutionary relationships of RNA virus pathogens identified in this study.**

506 Phylogenetic trees were estimated using the maximum likelihood method implemented in

507 PhyML. Sequences identified from this study were marked with red solid circle. For larger trees,

508 we only show lineages or sub-lineages which contain sequences identified in this study.

509



510

511 **Fig. 5. Comparison of prevalence levels between healthy and control groups.** Boxplots for  
 512 patient and control groups for each pathogen identified in this study, including RNA virus  
 513 (green), DNA virus (yellow), bacteria (orange) and fungus (blue). For clarity, only non-zero  
 514 abundance levels are reported.