

1 **Rapid feedback on hospital onset SARS-CoV-2 infections** 2 **combining epidemiological and sequencing data**

3
4 Oliver T Stirrup*¹,
5 Joseph Hughes²,
6 Matthew Parker^{3,4,5},
7 David G Partridge^{6,7},
8 James G Shepherd²,
9 James Blackstone⁸,
10 Francesc Coll⁹,
11 Alexander J Keeley^{6,7},
12 Benjamin B Lindsey^{6,7},
13 Aleksandra Marek¹⁰,
14 Christine Peters¹⁰,
15 Joshua B Singer²,
16 The COVID-19 Genomics UK (COG-UK) consortium^{11**},
17 Asif Tamuri¹²,
18 Thushan I de Silva^{6,7},
19 Emma C Thomson^{2,13,14},
20 Judith Breuer*¹⁵

21
22 1. Institute for Global Health, University College London, London, UK
23 2. MRC-University of Glasgow Centre for Virus Research, Glasgow, Scotland, United Kingdom
24 3. Sheffield Bioinformatics Core, The University of Sheffield, Sheffield, UK
25 4. Sheffield Institute for Translational Neuroscience, The University of Sheffield, Sheffield, UK
26 5. Sheffield Biomedical Research Centre, The University of Sheffield, Sheffield, UK
27 6. Sheffield Teaching Hospitals NHS Foundation Trust, Sheffield, UK
28 7. The Florey Institute for Host-Pathogen Interactions & Department of Infection, Immunity and
29 Cardiovascular Disease, Medical School, University of Sheffield, Sheffield, UK
30 8. The Comprehensive Clinical Trials Unit at UCL, University College London, London, UK
31 9. Department of Infection Biology, Faculty of Infectious and Tropical Diseases, London School
32 of Hygiene & Tropical Medicine, London, UK
33 10. Clinical Microbiology, NHS Greater Glasgow and Clyde, Glasgow, UK
34 11. <https://www.cogconsortium.uk>
35 12. Research Computing, University College London, London, UK
36 13. Institute of Infection, Immunity and Inflammation, College of Medical, Veterinary and Life
37 Sciences, University of Glasgow, Glasgow, UK
38 14. Department of Infectious Diseases, Queen Elizabeth University Hospital, Glasgow, UK
39 15. Division of Infection and Immunity, University College London, London, UK

40
41 *authors for correspondence (email: oliver.stirrup@ucl.ac.uk; j.breuer@ucl.ac.uk).

42 **Full list of consortium names and affiliations are in the Appendix.

43

Abstract

Background: Rapid identification and investigation of healthcare-associated infections (HCAIs) is important for suppression of SARS-CoV-2, but the infection source for hospital onset COVID-19 infections (HOCl) cannot always be readily identified based only on epidemiological data. Viral sequencing data provides additional information regarding potential transmission clusters, but the low mutation rate of SARS-CoV-2 can make interpretation using standard phylogenetic methods difficult.

Methods: We developed a novel statistical method and sequence reporting tool (SRT) that combines epidemiological and sequence data in order to provide a rapid assessment of the probability of HCAI among HOCl cases (defined as first positive test >48 hours following admission) and to identify infections that could plausibly constitute outbreak events. The method is designed for prospective use, but was validated using retrospective datasets from hospitals in Glasgow and Sheffield collected February-May 2020.

Results: We analysed data from 326 HOCl. Among HOCl with time-from-admission ≥ 8 days the SRT algorithm identified close sequence matches from the same ward for 160/244 (65.6%) and in the remainder 68/84 (81.0%) had at least one similar sequence elsewhere in the hospital, resulting in high estimated probabilities of within-ward and within-hospital transmission. For HOCl with time-from-admission 3-7 days, the SRT probability of healthcare acquisition was >0.5 in 33/82 (40.2%).

Conclusions: The methodology developed can provide rapid feedback on HOCl that could be useful for infection prevention and control teams, and warrants further prospective evaluation. The integration of epidemiological and sequence data is important given the low mutation rate of SARS-CoV-2 and its variable incubation period.

Abstract word count: 250

Funding: COG-UK HOCl funded by COG-UK consortium, supported by funding from UK Research & Innovation, National Institute of Health Research and Wellcome Sanger Institute.

Keywords: COVID-19; healthcare associated; hospital onset; nosocomial; outbreak; SARS-CoV-2; whole genome sequencing.

81 Introduction

82 Nosocomial transmission of SARS-CoV-2 presents a significant health risk to both vulnerable
83 patients and to healthcare workers (HCWs)^[1-5]. There is a variable incubation period, extending
84 up to day 14 from exposure to the virus in symptomatic cases^[6]. It is also known that transmission
85 is possible from asymptomatic or presymptomatic carriers^[7-10], complicating identification of
86 hospital-acquisition among hospital onset COVID-19 infections (HOCl) and tracing of likely
87 sources of infection.

88
89 There is now substantial evidence from retrospective studies that genome sequencing of
90 epidemic viruses, together with standard infection prevention and control (IPC) practice, better
91 excludes nosocomial transmissions and better identifies routes of transmission than IPC
92 investigation alone^[11-13]. The development of rapid sequencing methods capable of generating
93 pathogen genomes within 24-48 hours has recently created the potential for clinical IPC decisions
94 to be informed by genetic data in near-real-time^[14]. Although SARS-CoV-2 has a low mutation
95 rate^[15], sufficient viral diversity exists for viral sequences to provide information regarding potential
96 transmission clusters^[16]. However, phylogenetic methods alone cannot reliably identify linked
97 infections, and the need for clinical teams to gather additional patient data presents challenges to
98 the timely interpretation of SARS-CoV-2 sequence data.

99
100 To overcome these barriers, we have developed a sequence reporting tool (SRT) that integrates
101 genomic and epidemiological data from HOCl) to rapidly identify closely matched sequences
102 within the hospital and assign a probability estimate for nosocomial infection. The output report is
103 designed for prospective use to reduce the delay from sequencing to impact on IPC practice. The
104 work was conducted as part of the COVID-19 Genomics (COG) UK initiative, which sequences
105 large numbers of SARS-CoV-2 viruses from hospitals and the community across the UK^[17]. Here
106 we describe the performance of the SRT using COG-UK sequence data for HOCl cases collected
107 from Glasgow and Sheffield between February and May 2020 and explore how it may have
108 provided additional useful information for IPC investigations.

109

110 Methods

111 The SRT methodology is applied to HOCl cases, defined here as inpatients with first positive
112 SARS-CoV-2 test or symptom onset >48 hours after admission, without suspicion of COVID-19
113 at admission. The SRT algorithm returns an estimate of the probability that each HOCl acquired
114 their infection post-admission within the hospital, with information provided on closely matching
115 viral sequences from the ward location at sampling and wider hospital. Results for individual
116 HOCl are evaluated in relation to the IPC classification system recommended by Public Health
117 England (PHE), based on interval from admission to positive test: 3-7 days post admission =
118 indeterminate healthcare-associated infection (HCAI); 8-14 days post admission = probable
119 HCAI; >14 days post admission = definite HCAI^[18]. We also applied the PHE definition of
120 healthcare-associated COVID-19 outbreaks^[18] (i.e. ≥ 2 cases associated with specific ward, with
121 at least one being a probable or definite HCAI) to ward-level data, and for each outbreak
122 evaluated whether there was one or more distinct genetic cluster. This was determined by
123 consecutive linkage of each HOCl into clusters using a 2 SNP threshold (with HOCl assigned
124 to a genetic cluster if a sequence match to any member). Sequences with <90% genomic
125 coverage were excluded from all analyses.

126

127 Research Ethics

128 Research Ethics for COG-UK Consortium and research undertaken under its auspices was
129 granted by the PHE Research Ethics and Governance group as part of the emergency response
130 to COVID-19 (24 April 2020, REF: R&D NR0195) and by the relevant Scottish biorepository
131 authorities (16/WS/0207NHS and 10/S1402/33). This was a retrospective analysis on fully
132 anonymized data, the collection of which did not involve any active research intervention.
133 Consent therefore was neither required nor requested from individual patients.

134

135

136 Data collection and processing

137

138 *Glasgow*

139 During the first wave of SARS-CoV-2, the MRC-University of Glasgow Centre for Virus
140 Research collected residual clinical samples from SARS-CoV-2 infected individuals following
141 diagnosis at the West of Scotland Specialist Virology Centre. Samples were triaged for rapid
142 sequencing using Oxford Nanopore Technologies (ONT) for suspected healthcare related
143 infections or Illumina sequencing in all other cases (details in Appendix).

144

145 *Sheffield*

146 Residual clinical samples from SARS-CoV-2 positive cases diagnosed at Sheffield Teaching
147 Hospitals NHS Foundation Trust were sequenced at the University of Sheffield using ARTIC
148 network protocol^[19] and ONT. Throughout the epidemic, members of the IPC team were notified
149 by the laboratory and by clinical teams of positive results and reviewed relevant areas to ensure
150 optimisation of practice and appropriate management of patients. Electronic reports were
151 created contemporaneously, including an assessment as to whether suspected linked cases
152 were present based on ward level epidemiology. As part of SRT validation, these reports were
153 accessed retrospectively by a study team member blind to the sequencing data and each

154 included HOCl case was defined as being thought unlinked to other cases, a presumed index
 155 case in an outbreak or a presumed secondary case.

156

157

158 HOCl classification algorithm

159 The sequence matching and probability score algorithm is run separately for each ‘focus
 160 sequence’ corresponding to a HOCl. We use associated metadata to assign other previously
 161 collected sequences to categories representing where the individual may be part of a SARS-
 162 COV-2 transmission network:

- 163 • Unit reference set: individual could be involved with transmission on same unit
 164 (ward/ICU etc) as focus sequence (look-back interval: 3 weeks)
- 165 • Institution reference set: individual could be involved with transmission in same
 166 institution/hospital as focus sequence (look-back interval: 3 weeks)
- 167 • Community reference set: individual could be involved with transmission outside of focus
 168 sequence institution (look-back interval: 6 weeks).

169 It is possible for samples to be members of multiple reference sets. For example an outpatient
 170 may be involved in SARS-CoV-2 transmission at the institution they attended and/or in
 171 community transmission.

172

173 For each run of the algorithm, pairwise comparisons are conducted between the focus
 174 sequence and each sequence within the unit reference set, institution reference set and
 175 community reference set. A reference set sequence is considered a close match to the focus
 176 sequence if there is a maximum of two SNP differences between them. This choice was based
 177 on reported healthcare-associated outbreak events^[14, 20] and the overall mutation rate of SARS-
 178 CoV-2 (details in Appendix).

179

180 *Probability calculations*

181 We use an expression of Bayes theorem to estimate probabilities for post-admission infection of
 182 each focus case divided by exposure on the unit, within the rest of the institution and from
 183 visitors (if allowed). An estimate of the prior probability (P_{prior}) of post-admission infection for
 184 each focus case is modified to a posterior probability according to information provided by the
 185 sequence data. The algorithm is based on sound statistical principles, but involves heuristic
 186 approximations.

187

188 In symptomatic focus cases we base P_{prior} on the time interval (t) from admission to date of
 189 symptom onset or first positive test (if date of symptom onset not recorded). We calculate $P_{prior} =$
 190 $F(t)$, where $F(t)$ is the cumulative distribution function of incubation times^[6] (derivation in
 191 Appendix).

192

193 In theory, it would be optimal to use all of the information in the *exact* sequences observed.
 194 However, with the goal of constructing a computationally simple algorithm, we base our
 195 calculations on the probability of observing a *similar* sequence (within 2 SNPs) to that actually
 196 observed for each focus case conditional on each potential infection source/location: infection in
 197 the community, current unit/ward or elsewhere in the hospital/institution, or from a visitor. For

198 the unit and hospital, we estimate this probability using the observed sequence match
199 proportion (on pairwise comparison to the focus sequence) in the unit reference set and
200 institution reference set, respectively. For community- or visitor-acquired infection we use a
201 weighted proportion of matching sequences in the community reference set, with weightings
202 determined by a calibration model that describes geographic clustering of similar sequences
203 among community-acquired infections (described in Appendix). The geographic weighting
204 model was fitted separately for each study site using sequences strongly thought to represent
205 community-acquired infection: all community-sampled sequences and patients presenting to the
206 Emergency Department with COVID-19, excluding those recorded as being healthcare workers.
207

208 Software

209 The analysis was conducted in R (v. 4.0.2, R Foundation, Vienna), using sequence processing
210 and comparison functions from *ape* (v5.4) and geospatial functions in the *PostcodesioR* (v0.1.1)
211 and *gmt* packages (v2,0). R code to run the algorithm is available^[21], and it has also been
212 implemented as a standalone SRT for prospective use^[22] within COV-GLUE^[23].
213
214
215

216 Results

217 Study populations

218

219 *Glasgow*

220 The Glasgow dataset included 1199 viral sequences (available as of 23rd June 2020): 426 were
221 derived from community sampling sites, 351 from patients presenting to Emergency Department
222 or acute medical units, 398 from hospital inpatients and 24 from outpatients. Limited data were
223 available regarding the total number of HCWs testing positive and their identification among
224 community samples, but 15 sequences were recorded as being from HCWs. First positive test
225 dates ranged from 3rd March to 27th May 2020. All consensus sequences had genomic
226 coverage >90%.

227

228 We applied the SRT algorithm to data from three hospitals with required metadata available, for
229 which 128/246 inpatient cases with sequences were HOCIs. Two of these patients had been
230 transferred from another hospital within 14 days prior to their positive test and were not
231 processed as focus sequences. One inpatient without recorded sampling location was excluded,
232 leaving 125 HOCIs for analysis. Population sequencing coverage was 536/1578 (34.0%) overall
233 for patients at the three hospitals and 128/328 (39.0%) for HOCIs specifically (Appendix-figure
234 1).

235

236

237 *Sheffield*

238 The Sheffield dataset included 1630 viral sequences with accompanying metadata (available as
239 of 10th October 2020): 714 were from inpatients, 117 were from outpatients and 799 were from
240 HCWs. For this retrospective evaluation, 447/714 inpatient samples taken on date of admission
241 were assumed to represent community-onset cases and used to calibrate the model. First
242 positive test dates ranged from 23rd February to 30th May 2020. One sequence with genome
243 coverage <90% was dropped from further analysis (an inpatient on date of admission). 201 of
244 the inpatients were HOCIs. Population sequencing coverage was 714/977 (73.1%) overall for
245 inpatients, 201/261 (77.0%) for HOCIs specifically and 799/962 (83.1%) for HCWs.

246

247 Comparison to standard PHE classification

248 SRT algorithm results in comparison to standard PHE classifications are summarised in Figure
249 1 and Table 1. The majority of HOCI cases in Glasgow (78/125, 62.4%) and over a third in
250 Sheffield (71/201, 35.3%) met the definition of a definite HCAI and so are known to have
251 acquired the virus post-admission irrespective of sequencing results. The probable HCAI cases
252 formed the next largest group at each site. Overall, the SRT algorithm identified close sequence
253 matches from the same ward for 66.4% of definite and 64.2% of probable HCAIs, indicating
254 likely within-ward transmission (examples in Case Studies). When one or more close sequence
255 match was identified on the focus sequence's ward, the SRT probability of infection on the ward
256 was >0.5 in 185/189 cases (Figure 2). For indeterminate HCAIs the SRT probability of HCAI
257 was >0.5 in 33/82 (40.2%), and in 27/33 (81.8%) a close sequence match on the ward was
258 present. Overall, 14/125 (11.2%) HOCIs in Glasgow and 175/201 (87.1%) in Sheffield had at

259 least one close sequence match to a HCW sample, reflecting the much greater availability of
260 sequences from HCWs in the Sheffield dataset.

261
262 In 16/244 (6.6%) cases that met the probable or definite HCAI definitions, there was no
263 sequence match within the hospital; this is likely due to incomplete sequence data from SARS-
264 CoV-2 hospitalised cases and staff (with population sequencing coverage <40% patients and
265 very limited for staff from Glasgow and ≈75% of patients and staff in Sheffield) and the presence
266 of asymptomatic and/or undiagnosed carriers. To reflect this the SRT will report “This is a
267 probable/definite HCAI based on admission date, but we have not found genetic evidence of
268 transmission within the hospital” in such situations. There were 26 HOCl in the Sheffield
269 dataset for whom it was recorded that visitors were allowed on the ward at time of sampling. In
270 three of these the estimated probability of infection from a visitor was between 0.4 and 0.5 (all
271 had ≥18 days from admission and no ward close sequence matches).

272
273 Within the Sheffield dataset we identified six wards with two genetically distinct outbreak
274 clusters (of two or more patients) and three wards with three distinct outbreaks (see Case Study
275 2). Standard IPC assessment had classified each as a single outbreak. We also identified 10
276 and 44 HOCl in the Glasgow and Sheffield datasets, respectively, with no apparent genetic
277 linkage to other HOCl cases on the ward but who met the PHE definition of inclusion within an
278 outbreak event (Table 2).

279

280 Comparison to local IPC conclusions in Sheffield

281 Contemporaneous notes by IPC teams in Sheffield classified 18/201 HOCIs as the index case
282 in outbreaks. IPC staff defined an index case as the first detected in an environment regardless
283 of prior inpatient stay and, correspondingly, of these 14/18 were the first sequence on their ward
284 and one was the second (the first 1 day earlier from a different bay on the ward was also
285 recorded as an index case, and IPC staff deemed a ward outbreak with unclear index or
286 possibly 2 index cases). Of the 18 index cases 11 showed at least one subsequent close
287 sequence match on the same ward (the 2 index cases on a single ward were not genetically
288 similar, and for 1/18 there were no subsequent sequences from the ward). The median SRT
289 probability of HCAI was 0.70 (IQR 0.22-1.00, range 0.04-1, >0.5 in 12/18).

290

291 A further 144/201 HOCIs were classified as being part of local outbreaks, and among these the
292 median SRT probability of HCAI was 0.98 (IQR 0.89-1.00; range 0.02-1.00; >0.5 in 129/144)
293 with one or more close sequence match on the same ward in 104/144. The remaining 39/201
294 HOCIs, including 10 that were not recorded as HOCIs at the time, were classified by the IPC
295 teams as not being part of local outbreaks. Among these the median SRT probability of HCAI
296 was 0.74 (IQR 0.23-0.99, range 0.02-1.00; >0.5 in 23/39), with one or more close sequence
297 matches on the same ward in 7/39.

298

299 *Case Study 1*

300 Figure 3 shows a phylogenetic tree of eight HOCIs within a single ward at a Glasgow hospital
301 (Hospital 5, Unit 93), alongside associated meta-data and SRT probability outputs. The first
302 HOCI detected (UID0032) was transferred from another hospital within the previous 2 weeks
303 and so SRT output was not generated. All subsequent HOCIs return close sequence matches to
304 at least one prior case on the ward, leading to SRT probability estimates of ward-acquired
305 infection >0.9, even for UID0017 (an indeterminate HCAI). The phylogenetic tree indicates
306 UID0032 has a SNP lacked by most of the cases identified on the ward, and therefore did not
307 seed all of the cases in the outbreak cluster. Also shown is a single HOCI from a different ward
308 in the same hospital (UID0025); this individual was an indeterminate HCAI, but a higher
309 proportion of similar viral sequences within the hospital in comparison to their local community
310 led to a SRT result of probable hospital-acquired infection.

311

312 *Case Study 2*

313 Figure 4 shows phylogenetic trees relating to three distinct viral lineages identified on a single
314 ward in the Sheffield dataset (classified by contemporaneous IPC investigation as a single
315 outbreak). Two of these lineages also include sequences from inpatients sampled from other
316 wards within the same hospital. Detailed ward movement data highlighted additional possible
317 links between patients in the B.2.1 cluster. Both UID0149 and UID0157 were present at
318 LOC0111 prior to their sample dates.

319 Discussion

320 We have developed a novel approach for identification and investigation of hospital-acquired
321 SARS-CoV-2 infections combining epidemiological and sequencing data, designed to provide
322 rapid and concise feedback to IPC teams working to prevent nosocomial transmission. Through
323 retrospective application to clinical datasets, we have demonstrated that the methodology is
324 able to provide confirmatory evidence for most PHE-defined definite and probable HCAs and
325 provide further information regarding indeterminate HCAs. Thus the SRT may allow IPC teams
326 to optimise their use of resources on areas with likely nosocomial acquisition events.

327
328 While the SRT is not likely to change IPC conclusions in cases meeting the definition of 'definite'
329 or 'probable' HCAI based on interval from admission to symptom onset, in 91% of cases it did
330 identify patients in the same ward or elsewhere in the hospital who could plausibly be linked to
331 the HOI within a single outbreak event. Those definite and probable HOIs without close
332 sequence matches are likely to reflect transmission from sources within the hospital that have
333 either not been diagnosed or who were diagnosed without viral sequencing. In such cases it is
334 impossible to calculate a probability of transmission and the SRT will simply state that no
335 sequence matches were found within the hospital.

336
337 For cases meeting the definition of 'indeterminate healthcare associated', the probability scores
338 returned would be useful for IPC teams. These probabilities are dependent on comparison to
339 sequences from cases of community-acquired infection obtained either from direct community
340 sampling or from patients sampled at admission. The Sheffield dataset was lacking the former
341 data source, but the SRT nonetheless classified a similar proportion of 'indeterminate
342 healthcare associated' HOIs as community-acquired infections to that found in the Glasgow
343 dataset (approximately 60%).

344
345 Current PHE guidelines define healthcare-associated COVID-19 outbreaks as two or more
346 cases associated with a specific setting (e.g. ward), with at least one case having illness onset
347 after 8 days of admission^[18]. However, the guidelines note that "investigations of healthcare
348 associated SARS-CoV-2 infection should also take into account COVID-19 cases categorised
349 as 'indeterminate healthcare associated'" (i.e. onset 3-7 days after admission), for which our
350 SRT output would be useful. In most HOIs meeting this definition of inclusion within an
351 outbreak event, we found evidence of clusters of similar viral sequences located on the ward
352 concerned, and the SRT results were in line with available local IPC classifications in the
353 majority of cases. However, a substantial minority (54/279) of HOIs although assumed to be
354 part of a ward outbreak, were, in fact, isolated cases for which the sequencing data refuted
355 genetic linkage to other sequences from the ward. The SRT also provided evidence of wards
356 where IPC-defined outbreak events comprised two or three clearly distinct viral lineages.

357
358 The retrospective datasets analysed in this study represent the first few months of the COVID-
359 19 epidemic in the UK, and nosocomial transmission of the virus in the UK during this period
360 has previously been reported at multiple sites^[14, 24, 25]. HCWs were at increased risk of infection
361 and adverse health outcomes^[1, 2, 4, 5, 26] and could have been important drivers of nosocomial
362 transmission^[8]. Data were limited for Glasgow but the Sheffield dataset contained a large

363 number of sequences obtained from HCWs, with population sequencing coverage for this group
364 >80%, and there was a close sequence match to at least one HCW observed for 87% of HOCl.
365 Our analysis has not evaluated direction of transmission to or from HCWs, but they were clearly
366 linked into transmission networks within the hospital. A limitation of the current SRT approach
367 and of the retrospective data available is that they do not include detailed information regarding
368 work locations for HCWs. However, prospective use of the SRT would allow IPC teams to
369 investigate linkage from a HOCl to any HCWs flagged as having a close sequence match.

370
371 While a phylogenetic approach is useful in excluding direct transmission between cases, it can
372 be more problematic to confirm transmission source^[27]. Phylogenetic models can evaluate the
373 full genetic information provided by viral sequence data, but there are challenges in
374 incorporating and summarising associated patient meta-data in a timely fashion^[28]. The
375 challenge of timely collection and standardisation of patient meta-data is also relevant for use of
376 the SRT that we have developed, but it is possible to automate such processes through
377 electronic patient record systems. There have been advances in recent years in the
378 computational efficiency and workflow standardisation possible for phylogenetic analyses that
379 have made it easier to use these methods for real-time investigation of outbreaks, for example
380 through the development of the Nextstrain project^[29, 30]. However, there does not currently exist
381 phylogenetic software for SARS-CoV-2 that produces reports or other outputs designed for
382 direct and immediate use by IPC professionals. There will be cases in which phylogenetic
383 analysis would provide information beyond that returned by the SRT, and the two approaches
384 may be complementary to one another for outbreak investigation.

385
386 Comparison of SRT output to phylogenetic trees in a number of test cases suggested that some
387 clusters of genetically similar cases identified within a specific ward likely represented more than
388 one transmission event onto the ward from similar viral lineages circulating within the healthcare
389 system. Whilst monophyletic clusters associated with a single location are easier to interpret, we
390 consider the presence of viruses within a ward or hospital that are genetically similar to a HOCl
391 as evidence for nosocomial infection even when they are not plausible transmission sources
392 themselves, given the potential for asymptomatic transmission^[7-10] and complex transmission
393 networks^[14].

394
395 The SRT uses a number of heuristic approximations in order to provide an integrated summary
396 of epidemiological and sequence data. However, this choice is associated with the limitation that
397 it does not provide a full probabilistic model of potential transmission networks. Further
398 development of the SRT would also aim to more fully incorporate patient movement data and
399 shift locations for HCWs.

400
401 We believe that collaboration between methodologists, virologists, IPC clinicians and software
402 engineers is essential in order to create workflows and reporting systems that will enable the
403 routine use of pathogen sequence data for IPC. The SRT represents such a collaboration, and it
404 has been designed to enable automation of the linkage and processing of viral sequence and
405 patient meta-data and subsequent feedback of relevant information to IPC staff. The automated
406 feedback provided by the SRT is nonetheless dependent on timely sequencing of a high

407 proportion of viral samples from cases within the hospital concerned, ideally in combination with
408 sequences also available from community-sampled cases. In the UK this has been possible
409 through the national COG-UK project^[17]. Denmark has also implemented high population-
410 coverage sequencing of SARS-CoV-2^[31], but this is not the case for most countries. The
411 emergence and rapid dominance of lineage B.1.1.7 in the UK^[32] has provided a case study for
412 the impact of national-level genomic surveillance, but further evidence is required to determine
413 whether rapid sequencing is worth the necessary investment for routine use within IPC practice.
414 This judgement would also be dependent on the available health infrastructure and resources at
415 both the local and national levels.

416
417 Prospective evaluation of the SRT is currently underway within a multicentre study in the UK^[33].
418 This study and its accompanying research program will evaluate the impact of routine viral
419 sequencing and use of the SRT on IPC knowledge, actions and outcomes, and will include
420 quantitative, qualitative^[34] and health economic analyses to help guide the future development
421 of pathogen genomics for IPC.

422
423 Our novel approach to the investigation of HOCl has shown promising characteristics on
424 retrospective application to two clinical datasets. The SRT described allows rapid feedback on
425 HOCl that integrates epidemiological and sequencing data to generate a simplified report at
426 the time that sequence data become available. Prospective evaluation is required in order to
427 recommend use of the SRT in clinical practice, and this work is ongoing. The methodology has
428 been developed for hospital inpatients, but the principles may also be applicable to other
429 settings.

430
431

432 Declaration of interests

433 JBr receives funding from NIHR, FC from Wellcome and MDP from NIHR. The remaining
434 authors do not have any declarations of interest.

435

436 Acknowledgements

437 COG-UK HOCl is funded by the COG-UK consortium, which is supported by funding from the
438 Medical Research Council (MRC) part of UK Research & Innovation (UKRI), the National
439 Institute of Health Research (NIHR) and Genome Research Limited, operating as the Wellcome
440 Sanger Institute. JBr receives funding from the NIHR ULC/UCLH Biomedical Research Centre.
441 FC is funded by Wellcome (grant number: 201344/Z/16/Z). MDP is funded by the NIHR
442 Sheffield Biomedical Research Centre (BRC - IS-BRC-1215-20017). We acknowledge the help
443 of the UCL Comprehensive Clinical Trials Unit. The authors wish to thank the NHS Greater
444 Glasgow and Clyde and Sheffield Teaching Hospitals NHS Foundation Trust infection
445 prevention and control teams for provision of data. The authors thank Michael Chapman for his
446 assistance in the development of this project.

447

448 Data sharing

449 The sequence data analysed are included within publicly available datasets
450 (<https://www.cogconsortium.uk/data/>), and a list of the relevant sequence identification codes is
451 provided (Supplementary File 1). Due to data governance restrictions related to individual
452 patient data linked to genetic sequences it is not possible to publicly share the associated meta-
453 data. Requests for access to the data can be made by submission of a research proposal to the
454 COG-UK Steering Committee (contact@cogconsortium.uk).
455
456

457 **References**

- 458 1. Kursumovic E, Lennane S, Cook TM. **Deaths in healthcare workers due to COVID-19: the**
459 **need for robust data and analysis.** *Anaesthesia* 2020; 75(8):989-992.
- 460 2. Wang D, Hu B, Hu C, Zhu F, Liu X, Zhang J, et al. **Clinical Characteristics of 138**
461 **Hospitalized Patients With 2019 Novel Coronavirus–Infected Pneumonia in Wuhan,**
462 **China.** *JAMA* 2020; 323(11):1061-1069.
- 463 3. Leclerc Q, Fuller N, Knight L, null n, Funk S, Knight G. **What settings have been linked to**
464 **SARS-CoV-2 transmission clusters? [version 2; peer review: 2 approved].** *Wellcome Open*
465 *Research* 2020; 5(83).
- 466 4. The DELVE Initiative. **Scoping Report on Hospital and Health Care Acquisition of**
467 **COVID-19 and its Control.** DELVE Report No. 3. Published 06 July 2020. In. [http://rs-](http://rs-delve.github.io/reports/2020/07/06/nosocomial-scoping-report.html)
468 [delve.github.io/reports/2020/07/06/nosocomial-scoping-report.html](http://rs-delve.github.io/reports/2020/07/06/nosocomial-scoping-report.html); 2020.
- 469 5. Shah ASV, Wood R, Gribben C, Caldwell D, Bishop J, Weir A, et al. **Risk of hospital**
470 **admission with coronavirus disease 2019 in healthcare workers and their households:**
471 **nationwide linkage cohort study.** *BMJ* 2020; 371:m3582.
- 472 6. Lauer SA, Grantz KH, Bi Q, Jones FK, Zheng Q, Meredith HR, et al. **The Incubation Period**
473 **of Coronavirus Disease 2019 (COVID-19) From Publicly Reported Confirmed Cases:**
474 **Estimation and Application.** *Annals of Internal Medicine* 2020; 172(9):577-582.
- 475 7. He X, Lau EHY, Wu P, Deng X, Wang J, Hao X, et al. **Temporal dynamics in viral**
476 **shedding and transmissibility of COVID-19.** *Nature Medicine* 2020; 26(5):672-675.
- 477 8. Rivett L, Sridhar S, Sparkes D, Routledge M, Jones NK, Forrest S, et al. **Screening of**
478 **healthcare workers for SARS-CoV-2 highlights the role of asymptomatic carriage in**
479 **COVID-19 transmission.** *eLife* 2020; 9:e58728.
- 480 9. Oran DP, Topol EJ. **Prevalence of Asymptomatic SARS-CoV-2 Infection.** *Annals of*
481 *Internal Medicine* 2020; 173(5):362-367.
- 482 10. Lucey M, Macori G, Mullane N, Sutton-Fitzpatrick U, Gonzalez G, Coughlan S, et al. **Whole-**
483 **genome Sequencing to Track Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-**
484 **CoV-2) Transmission in Nosocomial Outbreaks.** *Clinical Infectious Diseases* 2020.
- 485 11. Brown JR, Roy S, Shah D, Williams CA, Williams R, Dunn H, et al. **Norovirus**
486 **Transmission Dynamics in a Pediatric Hospital Using Full Genome Sequences.** *Clinical*
487 *Infectious Diseases* 2018; 68(2):222-228.
- 488 12. Houldcroft CJ, Roy S, Morfopoulou S, Margetts BK, Depledge DP, Cudini J, et al. **Use of**
489 **Whole-Genome Sequencing of Adenovirus in Immunocompromised Pediatric Patients to**
490 **Identify Nosocomial Transmission and Mixed-Genotype Infection.** *The Journal of Infectious*
491 *Diseases* 2018; 218(8):1261-1271.
- 492 13. Roy S, Hartley J, Dunn H, Williams R, Williams CA, Breuer J. **Whole-genome Sequencing**
493 **Provides Data for Stratifying Infection Prevention and Control Management of**
494 **Nosocomial Influenza A.** *Clinical Infectious Diseases* 2019; 69(10):1649-1656.
- 495 14. Meredith LW, Hamilton WL, Warne B, Houldcroft CJ, Hosmillo M, Jahun AS, et al. **Rapid**
496 **implementation of SARS-CoV-2 sequencing to investigate cases of health-care**
497 **associated COVID-19: a prospective genomic surveillance study.** *The Lancet Infectious*
498 *Diseases* 2020; 20(11):1263-1272.
- 499 15. Fauver JR, Petrone ME, Hodcroft EB, Shioda K, Ehrlich HY, Watts AG, et al. **Coast-to-**
500 **Coast Spread of SARS-CoV-2 during the Early Epidemic in the United States.** *Cell* 2020;
501 181(5):990-996.e995.
- 502 16. van Dorp L, Acman M, Richard D, Shaw LP, Ford CE, Ormond L, et al. **Emergence of**
503 **genomic diversity and recurrent mutations in SARS-CoV-2.** *Infection, Genetics and*
504 *Evolution* 2020; 83:104351.
- 505 17. The COVID-19 Genomics UK (COG-UK) consortium. **An integrated national scale SARS-**
506 **CoV-2 genomic surveillance network.** *The Lancet Microbe* 2020; 1(3):e99-e100.

- 507 18. Public Health England. **COVID-19: epidemiological definitions of outbreaks and**
508 **clusters in particular settings.** In: [https://www.gov.uk/government/publications/covid-19-](https://www.gov.uk/government/publications/covid-19-epidemiological-definitions-of-outbreaks-and-clusters)
509 [epidemiological-definitions-of-outbreaks-and-clusters](https://www.gov.uk/government/publications/covid-19-epidemiological-definitions-of-outbreaks-and-clusters); 2020.
- 510 19. ARTIC Network. **SARS-CoV-2 sequencing protocol.** In; 2020.
- 511 20. Rockett RJ, Arnott A, Lam C, Sadsad R, Timms V, Gray K-A, et al. **Revealing COVID-19**
512 **transmission in Australia by SARS-CoV-2 genome sequencing and agent-based**
513 **modeling.** *Nature Medicine* 2020; 26(9):1398-1404.
- 514 21. Stirrup O. **HOCI SRT R code software repository:**
515 https://github.com/ostirrup/HOCI_SRT_R_code. In; 2021.
- 516 22. HOCI Sequence Reporting Tool working group. **HOCI-COV-GLUE software repository:**
517 https://github.com/ostirrup/HOCI-COV-GLUE_15Nov2020. In; 2020.
- 518 23. Singer J, Gifford R, Cotten M, Robertson D. **CoV-GLUE: A Web Application for Tracking**
519 **SARS-CoV-2 Genomic Variation.** In: Preprints.org; 2020.
- 520 24. Rickman HM, Rampling T, Shaw K, Martinez-Garcia G, Hail L, Coen P, et al. **Nosocomial**
521 **Transmission of Coronavirus Disease 2019: A Retrospective Study of 66 Hospital-**
522 **acquired Cases in a London Teaching Hospital.** *Clinical Infectious Diseases* 2020.
- 523 25. Carter B, Collins JT, Barlow-Pay F, Rickard F, Bruce E, Verduri A, et al. **Nosocomial**
524 **COVID-19 infection: examining the risk of mortality. The COPE-Nosocomial Study (COVID**
525 **in Older PEople).** *Journal of Hospital Infection* 2020; 106(2):376-384.
- 526 26. Houlihan CF, Vora N, Byrne T, Lewer D, Kelly G, Heaney J, et al. **Pandemic peak SARS-**
527 **CoV-2 infection and seroconversion rates in London frontline health-care workers.** *The*
528 *Lancet* 2020; 396(10246):e6-e7.
- 529 27. Volz EM, Frost SDW. **Inferring the Source of Transmission with Phylogenetic Data.**
530 *PLOS Computational Biology* 2013; 9(12):e1003397.
- 531 28. Villabona-Arenas CJ, Hanage WP, Tully DC. **Phylogenetic interpretation during**
532 **outbreaks requires caution.** *Nature Microbiology* 2020; 5(7):876-877.
- 533 29. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. **Nextstrain: real-**
534 **time tracking of pathogen evolution.** *Bioinformatics* 2018; 34(23):4121-4123.
- 535 30. Huddleston J, Hadfield J, Sibley TR, Lee J, Fay K, Ilcisin M, et al. **Augur: a bioinformatics**
536 **toolkit for phylogenetic analyses of human pathogens.** *Journal of Open Source Software*
537 2021; 6(57):2906.
- 538 31. Bager P, Wohlfahrt J, Fonager J, Albertsen M, Yssing Michaelsen T, Holten Møller C, et al.
539 **Increased risk of hospitalisation associated with infection with SARS-CoV-2 lineage B.**
540 **1.1. 7 in Denmark.** 2021.
- 541 32. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, et al. **Assessing**
542 **transmissibility of SARS-CoV-2 lineage B.1.1.7 in England.** *Nature* 2021.
- 543 33. Blackstone J, Stirrup O, Mapp F, Panca M, Copas A, Flowers P, et al. **Protocol for the**
544 **COG-UK hospital onset COVID-19 infection (HOCI) multicentre interventional clinical**
545 **study: evaluating the efficacy of rapid genome sequencing of SARS-CoV-2 in limiting the**
546 **spread of COVID-19 in United Kingdom NHS hospitals.** *medRxiv*
547 2021:2021.2004.2013.21255342.
- 548 34. Flowers P, Mapp F, Blackstone J, Stirrup O, Breuer J. **Developing Initial Programme**
549 **Theory: COVID-19 Genomics UK Consortium Hospital-onset COVID-19 Study (COG-UK**
550 **HOCI).** DOI: 10.31235/osf.io/ysm35. *SocArXiv (pre-print)* 2021.

551

552 Tables

553 **Table 1** Summary of sequence reporting tool outputs for the Glasgow and Sheffield datasets,
 554 according to standard infection prevention and control (IPC) definitions recommended by Public
 555 Health England regarding likelihood of healthcare-associated infection (HCAI)

	Glasgow data			Sheffield data		
	IPC classification			IPC classification		
	Indeterminate HCAI	Probable HCAI	Definite HCAI	Indeterminate HCAI	Probable HCAI	Definite HCAI
<i>n</i> HOCl cases	20	27	78	62	68	71
Time from admission to sample*, days	4.5 (3-6)	11 (9-13)	48 (26-83)	5 (4-6)	9 (8-13)	22 (17-31)
<i>Summary of sequence matches returned for each HOCl case</i>						
Close sequence match on ward	5 (25.0)	15 (55.6)	53 (68.0)	24 (38.7)	46 (67.6)	46 (64.8)
No close sequence match on ward, but match within hospital	8 (40.0)	7 (25.9)	19 (24.4)	34 (54.8)	21 (30.9)	21 (29.6)
No close sequence match anywhere within hospital	7 (35.0)	5 (18.5)	6 (7.7)	4 (6.5)	1 (1.5)	4 (5.6)
Close sequence match to one or more HCW	1 (5.0)	0 (0)	13 (16.7)	55 (88.7)	61 (89.7)	59 (83.1)
No close sequence match anywhere within dataset	2 (10.0)	1 (3.7)	4 (5.1)	4 (6.5)	1 (1.5)	4 (5.6)
<i>Probability calculations</i>						
Prior probability of HCAI†	0.39 (0.11-0.66)	0.97 (0.92-0.99)	1.00 (1.00-1.00)	0.49 (0.29-0.66)	0.92 (0.86-0.99)	1.00 (1.00-1.00)
Posterior probability of HCAI‡	0.33 (0.02-0.67)	0.98 (0.96-1.00)	1.00 (1.00-1.00)	0.40 (0.11-0.80)	0.98 (0.93-1.00)	1.00 (0.99-1.00)
Posterior probability of HCAI* category						
Low (<30%)	10 (50.0)	4 (14.8)	2 (2.6)	25 (40.3)	0 (0)	0 (0)
Moderately low (≥30% & <50%)	2 (10.0)	0 (0)	0 (0)	12 (19.4)	0 (0)	0 (0)
Medium (≥50% & <70%)	4 (20.0)	0 (0)	0 (0)	4 (6.5)	5 (7.4)	3 (4.2)
High (≥70% & <85%)	3 (15.0)	0 (0)	0 (0)	8 (12.9)	7 (10.3)	2 (2.8)
Very high (≥85%)	1 (5.0)	23 (85.2)	76 (97.4)	13 (21.0)	56 (82.4)	66 (93.0)

556 Data shown as median (interquartile range) or *n* (%). *or first +ve test where known. †Based on
 557 time from admission. ‡From source on ward or within hospital. HOCl, hospital onset COVID-19
 558 infection; HCW, healthcare worker.

559 **Table 2** Summary of distinct outbreak events for the Glasgow and Sheffield datasets, according
 560 to standard Public Health England (PHE) definition and with the addition of sequence data

	Glasgow data	Sheffield data
<i>n</i> HOCl cases	125	201
<i>n</i> ward locations	44	38
<i>Sequence matches per HOCl case</i>		
<i>n</i> sequence matches from same ward, median (IQR, range)	1 (0-5, 0-12)	1 (0-4, 0-18)
<i>n</i> sequence matches from rest of hospital, median (IQR, range)	3 (1-8, 0-52)	27 (5-52, 0-150)
<i>Standard PHE definition of outbreak event</i>		
HOCl cases part of ward outbreak event, <i>n</i> (%)	95 (76.0)	184 (91.5)
<i>n</i> ward outbreak events	17	24
<i>n</i> HOCl cases per ward outbreak event, median (IQR, range)	4 (2-8, 2-17)	5 (3.5-10.5, 2-28)
Days from first to last case in outbreak, median (IQR, range)	8 (6-15, 0-31)	18 (13-34, 3-68)
<i>n</i> wards with more than one distinct outbreak event	0	0
<i>Outbreak events with sequence linkage</i>		
HOCl cases part of ward outbreak event, <i>n</i> (%)	85 (68.0)	140* (69.7)
<i>n</i> ward outbreak events	16	33
<i>n</i> HOCl cases per ward outbreak event, median (IQR, range)	3.5 (2-8, 2-16)	3 (2-4, 1-19)
Days from first to last case in outbreak, median (IQR, range)	6 (4-9, 0-15)	4 (2-8, 0-17)
<i>n</i> wards with more than one distinct outbreak event	0	9†

561 *Includes two HOCl cases which each showed a close sequence match to another case on the same
 562 ward with interval from admission to sample date ≤ 2 days. †In three wards there were three
 563 genetically distinct outbreak events. HOCl, hospital onset COVID-19 infection; IQR, interquartile
 564 range.

565 Figures

566 **Figure 1** Plot of the posterior probability of healthcare-associated infection (HCAI) for (a)
567 Glasgow and (b) Sheffield hospital onset COVID-19 infection cases from the sequence reporting
568 tool algorithm against the prior probability of HCAI based only on time from admission to
569 diagnosis, grouped by standard infection prevention and control classification recommended by
570 Public Health England. Marginal histograms are displayed with bin-widths of 0.05.

571
572
573 **Figure 2** Plot of the posterior probabilities of healthcare-associated infection (HCAI) estimated
574 using the sequence reporting tool algorithm from a source on the current ward versus a source
575 elsewhere in the hospital for (a) Glasgow and (b) Sheffield hospital onset COVID-19 infection
576 cases grouped by standard Public Health England classification. In cases where there are no
577 close sequence matches in the dataset (including among community cases), the results
578 returned are based solely on the priors and the metadata; this explains the fact that there are
579 some cases with estimated posterior probability of infection on the ward greater than 0.5 for
580 whom there were no sequence matches on the ward.

581
582
583 **Figure 3.** Maximum-likelihood phylogeny of the sequences found in Hospital 5 Unit 93 and Unit
584 92 up until the 16th of May of the Glasgow dataset. The black lines represent the time from
585 admission to sampling. The values below the line are the posterior probability for unit infection +
586 the posterior probability of hospital infection from the sequence reporting tool. The tip nodes are
587 coloured according to the local authority area of the community surveillance sequences (circles)
588 or of the patients (crosses).

589
590
591 **Figure 4.** Maximum-likelihood phylogeny of the sequences found in Location '0111' in the
592 Sheffield dataset, also including patients at several other ward locations. The tree tip nodes are
593 coloured according to ward locations. The black lines represent the time from admission to
594 sampling. The values below the line are the posterior probability for unit infection + the posterior
595 probability of hospital infection from the sequence reporting tool. The circle containing a number
596 represents community sequences that are identical and at the base of this lineage (n=36).

597

Appendix

Rapid feedback on hospital onset SARS-CoV-2 infections combining epidemiological and sequencing data, by Oliver T Stirrup, Joseph Hughes, Matthew Parker, David G Partridge, James G Shepherd, James Blackstone, Francesc Coll, Alexander J Keeley, Benjamin B Lindsey, Aleksandra Marek, Christine Peters, Joshua B Singer, The COVID-19 Genomics UK (COG-UK) consortium, Asif Tamuri, Thushan I de Silva, Emma C Thomson, Judith Breuer

Methods

Details of sequencing protocols

Glasgow

Sequencing with ONT followed the protocols developed by the ARTIC network (v1 and v2) <https://artic.network/ncov-2019>. The reads were aligned to the reference strain (MN908947) using minimap2 (<https://doi.org/10.1093/bioinformatics/bty191>) and denoised using nanopolish (<https://www.nature.com/articles/nmeth.3444>) prior to primer trimming and consensus calling with iVar using a minimum depth of 20 reads (<https://doi.org/10.1186/s13059-018-1618-7>). Sequencing with Illumina also used the ARTIC network protocol for amplicon generation but was followed by a DNA KAPA library preparation kit (Roche) and indexing with NEBNext multiplex oligos (NEB) using 7 PCR cycles. Libraries were pooled and loaded on a MiSeqV2 cartridge. Illumina reads were processed with the PrimalAlign pipeline (<https://github.com/rjorton/PrimalAlign>). Briefly, reads were trimmed using trim_galore (https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) aligned to the reference using BWA ([10.1093/bioinformatics/btp698](https://doi.org/10.1093/bioinformatics/btp698)). Then, amplicon primers were removed and the consensus called with a read depth of 10 using iVar (<https://doi.org/10.1186/s13059-018-1618-7>). Metadata associated with each sample was collated in a redcap database (<https://www.project-redcap.org/>).

Sheffield

Sequencing with ONT followed the protocols developed by the ARTIC network (v1 and v2) <https://artic.network/ncov-2019>. Following base calling, data were demultiplexed using ONT Guppy using a high accuracy model. Reads were filtered based on quality and length (400 to 700bp), then mapped to the Wuhan reference genome and primer sites trimmed. Reads were then downsampled to 200x coverage in each direction. Variants were called using nanopolish (<https://github.com/jts/nanopolish>) and used to determine changes from the reference. Consensus sequences were constructed using reference and variants called.

41 Further details of reference set definitions

42

43 *Data sources for algorithm*

44 There are two potential sources of data for the HOCI classification algorithm. Firstly, there are
45 *institution-sampled sequences*: these include all viral sequences from samples obtained within
46 the institution/hospital. These sequences are linked to meta-data providing basic information
47 regarding the patient concerned and details of the sample from which the sequence was
48 obtained. Secondly, there are *community-sampled sequences*: these include all relevant
49 sequences obtained from samples from testing within the local community. These sequences
50 are associated with a more limited set of linked meta-data describing date of sample, residential
51 outer postcode of subject and place of work if they are recorded as being a HCW.

52 *Unit reference set*

53 This data set comprises all institution sequences sampled on or ≤ 3 weeks prior to (or ≤ 2 days
54 after for the prospective version of the SRT) the sample date of the focus sequence and for
55 which both the institution and the unit is the same as that for the focus sequence.

56

57 *Institution reference set*

58 This data set comprises firstly all institution-sampled sequences from HCWs, outpatients and
59 inpatients diagnosed >48 h after admission for which the institution matches that of the focus
60 sequence sampled on or ≤ 3 weeks prior to (or ≤ 2 days after for the prospective version of the
61 SRT) the sample date of the focus sequence and for which the unit is either not the same as
62 that for the focus sequence or is missing. Secondly, the data set includes all institution-sampled
63 sequences from A&E patients or inpatients diagnosed ≤ 2 days after admission for which the
64 institutionID matches that of the focus sequence sampled between (inclusively) 3 weeks and 3
65 days prior to the sample date of the focus sequence and for which the unit is either not the
66 same as that for the focus sequence or is missing. Thirdly, this data set also includes the subset
67 of community-sampled sequences of healthcare workers at the same institution as the focus
68 sequence.

69

70 *Community reference set*

71 This data set comprises firstly all community-sampled sequences sampled on or ≤ 6 weeks prior
72 to (or ≤ 2 days after for the prospective version of the SRT) the sample date of the focus
73 sequence. This data set also includes institution-sampled sequences sampled on or ≤ 6 weeks
74 prior to (or ≤ 2 days after for the prospective version of the SRT) the sample date of the focus
75 sequence from all non-inpatient samples, and those inpatients for whom sample date and
76 symptom onset date (if recorded) are both ≤ 2 days after the admission date.

77

78

79 Note that some institution-sampled sequences will contribute to both the community reference
80 set and either the unit reference set or the institution reference set (e.g. outpatients sampled
81 within 3 weeks prior to the focus sequence would be included in both the community reference
82 set and the institution reference set). HCWs recorded among the community-sampled
83 sequences within ≤ 3 weeks prior to the sample date of the focus sequence will also be included

84 in both the community reference set and the institution reference set if their workplace matches
85 the institution of the focus sequence.

86

87 Formulae for probability calculations

88 *Posterior of unit-acquired infection (UI) =*

$$\frac{P_{prior} * P_u * P(seq \pm 2SNPs|UI)}{P_{prior} * P_u * P(seq \pm 2SNPs|UI) + P_{prior} * P_v * P(seq \pm 2SNPs|VI) + P_{prior} * (1 - P_u - P_v) * P(seq \pm 2SNPs|II) + (1 - P_{prior}) * P(seq \pm 2SNPs|CI)}$$

89

90 *Posterior of institution-acquired infection (II) =*

$$\frac{P_{prior} * (1 - P_u - P_v) * P(seq \pm 2SNPs|II)}{P_{prior} * P_u * P(seq \pm 2SNPs|UI) + P_{prior} * P_v * P(seq \pm 2SNPs|VI) + P_{prior} * (1 - P_u - P_v) * P(seq \pm 2SNPs|II) + (1 - P_{prior}) * P(seq \pm 2SNPs|CI)}$$

91

92 *Posterior of visitor-acquired infection (VI) =*

$$\frac{P_{prior} * P_v * P(seq \pm 2SNPs|VI)}{P_{prior} * P_u * P(seq \pm 2SNPs|UI) + P_{prior} * P_v * P(seq \pm 2SNPs|VI) + P_{prior} * (1 - P_u - P_v) * P(seq \pm 2SNPs|II) + (1 - P_{prior}) * P(seq \pm 2SNPs|CI)}$$

93

94 *Posterior of community-acquired infection (CI) =*

$$\frac{(1 - P_{prior}) * P(seq \pm 2SNPs|CI)}{P_{prior} * P_u * P(seq \pm 2SNPs|UI) + P_{prior} * P_v * P(seq \pm 2SNPs|VI) + P_{prior} * (1 - P_u - P_v) * P(seq \pm 2SNPs|II) + (1 - P_{prior}) * P(seq \pm 2SNPs|CI)}$$

95

96

97 With terms defined as follows,

98 P_{prior} : prior probability of post-admission infection for each focus case, based on time interval
99 from admission to date of symptom onset or first positive test

100 P_u : prior probability of UI given post-admission infection (set based on expert opinion)

101 P_v : prior probability of VI given post-admission infection (set based on expert opinion)

102 $P(seq \pm 2SNPs|infection\ source/location)$: probability of observing a *similar* sequence (within 2 SNPs) to
103 that actually observed for each focus case conditional on each potential infection
104 source/location (estimated from sequence reference sets)

105

106 When there is a close sequence match found in any of the defined reference sets, the posterior
107 probability estimates for UI, II, VI and CI will always sum to 1. However, when there is no close
108 sequence match in any of the reference sets the posterior probability calculations are not valid
109 and the algorithm will return the prior probabilities for each potential source/location of infection.

110

111 Further details regarding sequence matching process

112 The ± 2 SNP threshold for a close sequence match was initially based on reports of healthcare-
113 associated outbreak events for which this was the maximum pairwise difference within clusters
114 (Meredith: DOI:10.1101/2020.05.08.20095687 & Rockett: DOI: 10.1101/2020.04.19.048751).

115 The outbreak events described included sequences with up to around 3 weeks between first
116 and last samples. This SNP threshold is also supported by calculations using the overall
117 mutation rate of SARS-CoV-2. If we take the average mutation rate of the virus to be 24
118 SNPs/year (Nextstrain value 24th June, <https://nextstrain.org/ncov/global?l=clock>), then
119 assuming independent (Poisson distributed) mutation events, ignoring the chance of mutations
120 occurring at the same position in the genome and using a fixed generation time of 5 days then
121 there is an approximate:

122 72% chance of no new SNPs per generation
 123 24% chance of 1 new SNP per generation
 124 4% chance of 2 new SNPs per generation
 125 0.4% chance of 3 new SNPs per generation
 126

127 A 2 SNP threshold would therefore be expected to identify close sequence matches between
 128 direct transmission pairs in a large majority of cases. Ambiguous nucleotide positions will be
 129 considered to match if there is an overlap in the possible values for the two sequences. 'N'
 130 values recorded in either the focus sequence or comparison sequence will be considered to be
 131 a match at that position.
 132
 133

134 Further details of prior probability calculations for post-admission infection

135
 136 We calculate $P_{prior} = F(t)$, where $F()$ is the cumulative distribution function of a published log-
 137 normal distribution for incubation times (Lauer et al: doi:10.7326/M20-0504;
 138 $\mu=1.621$, $\sigma=0.418$). For symptomatic HOI cases, the IPC classifications recommended by PHE
 139 translate into the following value ranges for P_{prior} :

- 140 ● indeterminate HCAI: 0.11 (onset 3 days post-admission) to 0.78 (onset 7 days
 141 post-admission)
- 142 ● probable HCAI: 0.86 (onset 8 days post-admission) to 0.99 (onset 14 days post-
 143 admission)
- 144 ● definite HCAI: $P_{prior} \geq 0.995$
 145

146 For asymptomatic focus cases, we define our prior on the basis that some proportion of the
 147 cases detected will never become symptomatic (P_a) with the remainder going on to develop
 148 symptoms within the next few days ($1-P_a$). We then define our prior probability of post-admission
 149 infection in these cases as:
 150

$$P_{prior} = (1 - P_a) * F(t + c) + P_a * F(t)$$

151
 152 where t is the interval from admissionDate to sampleDate and c is a constant reflecting the
 153 average interval within which we expect symptoms to appear (among those cases in which they
 154 do). P_a is set at 0.4 based on the findings of a published review article (Oran and Topol:
 155 doi.org/10.7326/M20-3012), and c is set to 3 based on a combination of expert opinion of the
 156 study PIs, the known distribution of time from infection to symptom onset and expert experience
 157 of asymptomatic screening.
 158

159 *Source given post-admission infection*

160 The model requires prior values for the probability of UI and VI given post-admission infection:
 161 P_u and P_v , respectively. However, in specifying the model we define P_u' as the probability of UI
 162 given post-admission infection when there are no visitors allowed on the ward, in which case the
 163 probability of VI is zero and $P_v'=0$. If visitors are allowed on the ward for the focus case, then we
 164 set $P_v = P_u' \times (1 - P_v)$.

165
 166 Based on expert opinion of the clinical co-authors, P_u' is set to different values according to the
 167 unit/ward type of the focus sequence with single bed wards having a lower prior probability of
 168 unit post-admission infection than bay wards: 0.5 for single bed wards and 0.7 for bay wards.
 169 We assumed a P_v of 0.2. The P_u values (when visitors are allowed) are therefore: 0.4 for single
 170 bed wards and 0.56 for bay wards. The largest of the three Glasgow hospitals included
 171 comprises single-room wards, whilst the other two and the Sheffield site comprise bay wards.
 172

173 Derivation of prior probability for post-admission infection

174 If we assume a uniform individual-level hazard (λ) of infection from 1st February 2020 (t_0),
 175 whether in hospital or not, then the probability density function (PDF) of infection at time t_{inf} from
 176 this date is: $\lambda e^{-(\lambda t_{inf})}$. The PDF of infection at time t_{inf} conditional on this occurring at any point
 177 prior to the date of symptom onset (t_{onset}) is: $(\lambda e^{-(\lambda t_{inf})}) / (1 - e^{-(\lambda t_{pos})})$, which is approximately
 178 $1/t_{onset}$ for small λ (taking the limit as $\lambda \rightarrow 0$). For HOCl cases, we are interested in whether t_{inf}
 179 occurred before or after the time of admission to hospital (t_{adm}). Also considering the evidence
 180 provided by the known incubation time of the disease (PDF f and CDF F), we integrate over the
 181 range of possible infection dates:

$$\begin{aligned}
 182 \quad P(t_{adm} \leq t_{inf} \mid t_{inf} \leq t_{onset}, T_{onset} = t_{onset}) &= [\int_{t_{adm}}^{t_{onset}} f(t_{onset} - x)/t_{onset} \cdot dx] / [\int_0^{t_{onset}} f(t_{onset} - x)/t_{onset} \cdot dx] \\
 183 &\approx [\int_{t_{adm}}^{t_{onset}} f(t_{onset} - x)/t_{onset} \cdot dx] / [\int_{-\infty}^{t_{onset}} f(t_{onset} - x)/t_{onset} \cdot dx] \\
 184 &= [\int_{t_{onset}-t_{adm}}^0 -f(u)/t_{onset} \cdot du] / (1/t_{onset}) \\
 185 &= \int_0^{t_{onset}-t_{adm}} f(u) \cdot du \\
 186 &= F(t_{onset} - t_{adm}) \\
 187 \\
 188
 \end{aligned}$$

189 Geographic weighting for community reference set

190 *Geographic weighting function*

191 The weight of each sequence within the community reference set is determined by geographic
 192 distance from the residential outer postcode of the focus case, using a function of the form:
 193 $\text{weight} = (1 - \beta) \cdot \exp(-\tau \cdot \text{communityDistanceToIndex}[i]) + \beta$,
 194 where, β takes a value between 0 and 1, and $\tau > 0$. These parameters are set based on
 195 calibration to the available community reference set at each site. The rationale for this weighting
 196 is that there is likely to be geographic clustering of viral lineages, and so newly observed
 197 community transmissions of SARS-CoV-2 are more likely to show genetic similarity to past
 198 sequences from the local area of that individual's home than to past sequences from regions
 199 that are further away. If postcode is missing for a case in the community reference set, then
 200 distance to the focus sequence is set to 100 km.
 201

201

202

203 *Statistical model for derivation of geographic weighting parameters*

204 The statistical model for geographic weighting is fitted separately for each study site using
 205 sequences which are strongly thought to represent community-acquired infection: all
 206 community-sampled sequences and patients presenting to A&E with COVID-19, excluding

207 those who are recorded as being healthcare workers or who do not have an available valid
 208 outer postcode. We will refer to these sequences as the ‘calibration set’.

209
 210 A statistical model is constructed to find the optimal values of β and τ to maximise the estimated
 211 probability ($P_{sim:i}$) of a newly observed community-acquired case having a similar sequence
 212 (± 2 SNPs) to that observed for each sequence in the calibration set. The estimated probability in
 213 each case within the calibration set is calculated as a weighted sum of ‘close match’ indicator
 214 variables for all other sequences in the calibration set sample from 6 weeks prior up until the
 215 sample date of that case, with the weighting function defined in terms of geographic distance
 216 between residential outer postcodes and the β and τ parameters as described for the
 217 community reference set.

218
 219 An overall log-likelihood function is defined using a Bernoulli distribution for each of the n
 220 sequences within the calibration set:

$$221 \ell = \sum_{i=1}^n \log(P_{sim:i}).$$

222 The values of β and τ that maximise ℓ were obtained for each of the study sites using the
 223 ‘bbmle’ package for R, with logit-parameterisation of β and log-parameterisation of τ .

224
 225 We assume that the probability of a sequence match conditional on infection from visitor on
 226 unit/ward can be calculated using the same weighting scheme as for the probability of a
 227 sequence match conditional on community-acquired infection (i.e. $P(\text{seq}\pm 2$
 228 SNPs|CI) = $P(\text{seq}\pm 2 \text{ SNPs}|VI)$).

229

230 Additional matching on ward location history

231 There is the potential for the algorithm described to return large numbers of close sequences
 232 matches with the hospital as a whole, which may make it difficult for IPC teams to use the
 233 output to direct their investigations when there are no potential sources of infection identified on
 234 the same ward as the focus case. We propose a location matching procedure in order to
 235 highlight the most relevant sequence matches for further investigation. This process does not
 236 currently form part of the statistical model, meaning that it can be treated as optional
 237 functionality for the SRT in the COG-UK HOCl study, and we have restricted the input data to a
 238 simplified format in order to minimise data management requirements.

239

240 For each inpatient sample in the input meta-data for the algorithm, we specify a single string
 241 variable comprising the concatenated names of any ward locations in the ≤ 14 days prior to the
 242 sample date and a separate string variable with any ward locations in the ≤ 14 days after the
 243 sample date. For each focus case submitted to the algorithm, output is flagged if there is any
 244 match identified between the wards listed in each of these fields or the ward at time of sampling
 245 for a close sequence match in comparison to the prior and current ward locations for the focus
 246 sequence (excluding those cases where there is already matching ward location at time of
 247 sampling for each).

248

249

250 Details of phylogenetic methods

251 Phylogenies were produced by the grapevine pipeline (<https://github.com/COG-UK/grapevine>)
252 as part of the COG-UK Consortium (<https://www.cogconsortium.uk>). Briefly, sequences from
253 GISAID and those produced as part of the COG-UK Consortium are independently quality
254 controlled and aligned to the Wuhan reference using minimap2
255 (<https://doi.org/10.1093/bioinformatics/bty191>). The two alignments are then combined, the
256 homoplasy at site 11083 is masked and the tree is reconstructed using FastTreeMP
257 (<http://www.microbesonline.org/fasttree/>). For each of the hospitals of interest, the tree is pruned
258 to keep sequences from Scotland or Yorkshire (as relevant) and by date excluding sequences
259 subsequent to the last “focus” patient sample date on the ward.

260

261

262 Details of SRT report format

263 The SRT system for prospective use needs to provide useful and appropriate feedback in both
 264 low incidence and high incidence settings for new HOCl cases. This is planned through the
 265 generation of a concise one-page PDF summary report for each focus sequence. This summary
 266 report contains key focus sequence meta-data, information regarding the estimated probabilities
 267 for infection source and details of up to ten close sequence matches identified within the same
 268 unit/ward and/or elsewhere in the hospital.

269

270 Probability summary categories

271 The sequence matching and probability score algorithm generates probability estimates for the
 272 source of infection for the focus patient being from the current unit/ward, from elsewhere in the
 273 hospital, from the community (pre-admission) or from a visitor. These probability estimates
 274 always sum to 1. In the summary report, probability estimates for each source of infection are
 275 categorised using the following levels:

- 276 ● 0-30%: low
- 277 ● 30-50%: moderately low
- 278 ● 50-70%: probable
- 279 ● 70-85%: high
- 280 ● 85-100%: very high

281

282 For clarity of presentation and communication, probability categories will not always be
 283 displayed in the summary report for all four potential sources of infection (i.e. ward/unit,
 284 elsewhere in hospital, visitor, or community). Special handling rules for specific situations are
 285 described below.

286

287 Close sequence matches within the same unit and/or hospital

288 The maximum number of close sequence matches that can be listed on the one-page summary
 289 report is 10 (for the combined sum of unit-level and institution-level matches). If the number of
 290 ward-level matches is $n > 5$ and the total number of close sequence matches is $N > 10$, then the
 291 number of ward-level matches is truncated at $5 + \max((5 - (N - n)), 0)$. If there are over ten close
 292 sequence matches in total, then the following message is displayed "Over 10 close matches;
 293 see detailed report for further information".

294

295 Within each set of unit-level and institution-level close sequence matches, ordering and priority
 296 for inclusion within the available slots is determined by the following set of criteria (in decreasing
 297 order of importance):

- 298 1. Number of SNPs relative to Wuhan strain present in comparison sequence but absent in
 299 focus sequence (fewer = higher priority)
- 300 2. Number of SNPs relative to Wuhan strain present in focus sequence but absent in
 301 comparison sequence (fewer = higher priority)
- 302 3. Whether comparison sequence is from a HCW (HCWs listed first)
- 303 4. HCAI status of comparison sequence (priority order: definite, probable, indeterminate,
 304 otherwise)

- 305 5. Samples from the past before samples in future
 306 6. Samples from within the two weeks prior to focus sequence sample date before others
 307 7. Number of units overlapping with focus sample's units

308

309

310 Report messages for specific output combinations

311

312 *No close sequence matches on unit/ward*

313 If there are no close sequence matches to the focus sequence on their current unit/ward, then
 314 no probability category is reported for this potential infection source (the algorithm returns a zero
 315 probability in such cases, which could be misleading given uncertainty over screening and
 316 sequencing coverage). The message "No matches from within unit" is displayed. The probability
 317 score category for infection from elsewhere in the hospital is provided in such cases.

318

319 *No close sequence matches elsewhere in hospital*

320 If there are no close sequence matches to the focus sequence elsewhere in the hospital, then
 321 no probability category is reported for this potential infection source. The message "No matches
 322 elsewhere in hospital" is displayed.

323

324 *No evidence of transmission within unit or hospital for probable or definite HCAI*

325 If the estimated probability of community-acquired infection from the algorithm is >50%, but the
 326 interval from admission to symptom onset (if recorded) or sample date is ≥ 8 days, then the
 327 following message is displayed in place of the estimated probability of community-acquired
 328 infection "This is a probable/definite HCAI based on admission date, but we have not found
 329 genetic evidence of transmission within the hospital".

330

331 *Probable unit- or hospital-acquired infection with source unclear*

332 If the posterior probability of unit-acquired infection and the posterior probability of infection from
 333 a source elsewhere in the hospital are each estimated to be <50%, but the sum of these two
 334 posterior probabilities is $\geq 50\%$, then the following message is displayed "Overall, this is a
 335 probable unit- or institution-acquired infection with source unclear".

336

337 Timeline graph

338 The timeline graph provides a visual representation of available sequences from the same
 339 unit/ward and the same institution/hospital as the focus sequence in the period from 3 weeks
 340 prior to their sample date to 1 week after. The key indicates which sequences are close
 341 matches to the focus sequence, and the numbering corresponds to that in the tabular summary
 342 of most relevant close sequence matches.

343

344 Sequencing prioritisation for prospective use of the SRT

345 The SRT algorithm was initially designed for use with comprehensive sequencing of all SARS-
346 CoV-2 cases within a hospital, in combination with representative sequencing of community-
347 sampled cases. However, it may be difficult to achieve high population sequencing coverage in
348 some situations, such as if there is a sudden surge in new admissions to the hospital and/or in
349 new HOCl cases. In such scenarios, we have recommended the following prioritisation of
350 samples (from highest to lowest) for sequencing within the prospective HOCl study

351 (<https://clinicaltrials.gov/ct2/show/NCT04405934>):

- 352 1. HOCl cases
- 353 2. SARS-CoV-2 +ve patients on wards where there is a HOCl case
- 354 3. HCWs with known contact with HOCl cases
- 355 4. Other HCWs
- 356 5. SARS-CoV-2 +ve patients admitted to any other wards
- 357 6. SARS-CoV-2 +ve patients attending for acute care (e.g. Accident and Emergency) but
358 not admitted

359

360 These prioritisation rules are guided by the following rationale:

- 361 - Most probable and definite HCAs (based on time from admission) show a close sequence
362 match to at least one other case on the same ward, so sequencing of HOCl cases and any
363 cases on the same ward would be enough to identify these links.
- 364 - Links between ward outbreaks will be of particular importance to IPC investigations, and would
365 be identified with sequencing focused on HOCl cases.
- 366 - The probability calculations within the SRT are most important for indeterminate HCAs, and
367 where there is no sequence match on the same ward the estimated probability of nosocomial
368 infection is <50% in the majority of such cases (36/38 for the Sheffield dataset). The probability
369 estimates for indeterminate HCAs should be interpreted with caution where overall sequencing
370 coverage is poor, but SRT results are unlikely to lead to inappropriate changes to standard IPC
371 actions if groups '1', '2' and '3' have been sequenced.
- 372 - Where there is a complete lack of close sequence matches within the hospital for probable or
373 definite HCAs, the SRT returns the message that there is a lack of available genetic evidence
374 for linkage (but not that nosocomial infection is unlikely).

375

376 Following from this reasoning, we feel that useful information would be returned by the SRT as
377 long as high sequencing coverage is achieved for groups '1', '2' and '3'. High sequencing
378 coverage of groups '4', '5' and '6' would allow the SRT to identify potential links between cases
379 that would likely be missed by standard IPC investigations.

380

381 For indeterminate HCAs with no close sequence matches on the same ward, an inaccurate
382 'zero' posterior probability of post-admission infection will be returned if one or more similar
383 sequence is found in the community reference set but no similar sequences are observed in the
384 institution reference set with imperfect sequencing coverage. This is likely to be a more
385 important issue in the setting of low SARS-CoV-2 incidence.

386

387 For example, if there are 40 cases that could be included in the institution reference set for a
388 focus sequence and 2 of these (5%) would be a close sequence match, then we would need to
389 sequence at least 31/40 (77.5%) in order to have $\geq 95\%$ probability of observing at least one of
390 the close sequence matches. However, if there are 200 cases that could be included in the
391 institution reference set and 10 of these (5%) would be a close sequence match, then we would
392 need to sequence at least 51/200 (25.5%) in order to have $\geq 95\%$ probability of observing at
393 least one of the close sequence matches.

394
395 A similar relationship would also be observed if we consider a rarer sequence type. If there are
396 40 cases that could be included in the institution reference set for a focus sequence and 1 of
397 these (2.5%) would be a close sequence match, then we would need to sequence at least 38/40
398 (95%) in order to have $\geq 95\%$ probability of observing the one close sequence match. However,
399 if there are 200 cases that could be included in the institution reference set and 5 of these
400 (2.5%) would be a close sequence match, then we would need to sequence at least 90/200
401 (45%) in order to have $\geq 95\%$ probability of observing at least one of the close sequence
402 matches.

403
404 On this basis, we believe that the goal of close to 100% sequencing coverage should be
405 pursued in the setting of low incidence of SARS-CoV-2, but that overall sequencing coverage of
406 50% or more may be sufficient in the event that a high incidence of SARS-CoV-2 leads to too
407 great a case load for available sequencing resources.

408
409

410 Results

411

412 Sequencing coverage in Glasgow dataset

413

414 **Appendix-figure 1** Proportion of cases sequenced in Greater Glasgow and Clyde Health Board
415 between 1 March and 27th May (with sequence available as of 23 June 2020) by location of test
416 (A). Also displayed are the proportion of sequenced cases in the three focus hospitals
417 subdivided by assessment and inpatient locations (B), and the proportion of HOCI cases
418 sequenced at these hospitals (C).

419

420

421 Home residence locations and geographic model parameters

422

423 **Appendix-figure 2** Home residence location of individuals in (a) the Glasgow dataset and (b)
424 the Sheffield dataset, displayed by sample source (not including HCWs). Locations are
425 analysed using only the outer postcode, and as such random jitter (within longitude and latitude
426 of 0.05) has been added to allow display without overlap of points. Plot created using ggmap for
427 R with map obtained from Stamen maps. For Glasgow 766 cases were included in the
428 calibration set with estimates of $\tau=0.15$ and $\beta=0.0$ for the geographic clustering model, whilst for
429 Sheffield 446 cases were included in the calibration set with resulting estimates of $\tau=0.84$ and
430 $\beta=0.16$.

431

432

433 SNP distance distributions

434 For the Glasgow sequence dataset as a whole the median pairwise SNP difference among all
435 sequences was 9, and there were 1.3%, 3.4%, 6.4% and 10.1% of pairwise comparisons with 0,
436 ≤ 1 , ≤ 2 and ≤ 3 SNP differences, respectively. For the Sheffield dataset as a whole the median
437 pairwise SNP difference among all sequences was 8, and there were 1.2%, 3.3%, 6.5% and
438 10.8% of pairwise comparisons with 0, ≤ 1 , ≤ 2 and ≤ 3 SNP differences, respectively.

439

440

441 **Appendix-figure 3** Frequency plot of all pairwise SNP differences among (a) all 1199
442 sequences in the Glasgow dataset and (b) all 1629 analysed sequences in the Sheffield
443 dataset.

444

445

446

447 Additional Case Study

448
449 Appendix-figure 4 shows a phylogenetic tree indicating complex transmission networks across
450 multiple hospitals in the Glasgow area (with SRT outputs for Hospitals 2 and 4). A monophyletic
451 cluster of HOICs can be seen in Hospital 2 Unit 48, with the first detected case identified by the
452 SRT as a hospital-acquired and the rest unit-acquired infections. A paraphyletic group of HOICs
453 was detected in Hospital 4 Unit 69. Patient 1 (UID0042) was screened for COVID in Unit 69 on
454 14.04.20 after developing a cough and oxygen requirement. The patient was moved from the
455 nightingale area to a single room on the ward on 14.04.20 and was confirmed positive on
456 15.04.20.

457
458 On 20.04.20 a second patient on Unit 69 (not sequenced) was screened after developing a
459 cough and pyrexia and confirmed positive on 21.04.2020. The patient was in a single room at
460 the time of symptom onset, however they had been in the main nightingale ward opposite
461 patient 1 for 5 days. At this point 13 asymptomatic contacts in Unit 69 were screened, and 8
462 (UID0043, UID0073, UID0041, UID0095, UID0116, UID0094, UID0083, UID0121) were positive.
463 These cases are all identified as hospital-acquired or unit-acquired infections and can be
464 grouped into a genetically similar cluster with a maximum pairwise distance of 2 SNPs between
465 each member and its nearest neighbour. However, this cluster clearly represents multiple
466 introductions of SARS-CoV-2 onto the ward.

467
468
469 **Appendix-figure 4.** Maximum-likelihood tree for sequences found in Hospital 2 Unit 48 and
470 Hospital 4 Unit 69 of the Glasgow dataset up until the 21st of April (inclusive). The circles with
471 numbers represent the number of community sequences that are identical and at the base of
472 each lineage (n=5, n=35, n=4). Tree tips with black circles represent further community
473 sequences. The black lines represent the time from admission to sampling. The values below
474 the line are the posterior probability for unit infection + the posterior probability of hospital
475 infection from the sequence reporting tool.

476
477

478 Examples of SRT reports

479 **Appendix-figure 5.** Example of sequence reporting tool output with estimated very highly
480 probable infection within unit.

481
482 **Appendix-figure 6.** Example of sequence reporting tool output with estimated probable
483 infection within hospital.

484

485 Sequence list for analysis

486 **Supplementary-file 1.** Comma separated value file containing a list of the COG-UK
487 identification codes for viral sequences included in the analysis.

488

- 489 The COVID-19 Genomics UK (COG-UK) consortium
490
491 **Funding acquisition, leadership, supervision, metadata curation, project administration,**
492 **samples, logistics, Sequencing, analysis, and Software and analysis tools:**
493 Dr Thomas R Connor PhD^{33, 34}, and Professor Nicholas J Loman PhD¹⁵.
494
495 **Leadership, supervision, sequencing, analysis, funding acquisition, metadata curation,**
496 **project administration, samples, logistics, and visualisation:**
497 Dr Samuel C Robson Ph.D⁶⁸.
498
499 **Leadership, supervision, project administration, visualisation, samples, logistics,**
500 **metadata curation and software and analysis tools:**
501 Dr Tanya Golubchik PhD²⁷.
502
503 **Leadership, supervision, metadata curation, project administration, samples, logistics**
504 **sequencing and analysis:**
505 Dr M. Estee Torok FRCP^{8, 10}.
506
507 **Project administration, metadata curation, samples, logistics, sequencing, analysis, and**
508 **software and analysis tools:**
509 Dr William L Hamilton PhD^{8, 10}.
510
511 **Leadership, supervision, samples logistics, project administration, funding acquisition**
512 **sequencing and analysis:**
513 Dr David Bonsall PhD²⁷.
514
515 **Leadership and supervision, sequencing, analysis, funding acquisition, visualisation and**
516 **software and analysis tools:**
517 Dr Ali R Awan PhD⁷⁴.
518
519 **Leadership and supervision, funding acquisition, sequencing, analysis, metadata**
520 **curation, samples and logistics:**
521 Dr Sally Corden PhD³³.
522
523 **Leadership supervision, sequencing analysis, samples, logistics, and metadata curation:**
524 Professor Ian Goodfellow PhD¹¹.
525
526 **Leadership, supervision, sequencing, analysis, samples, logistics, and Project**
527 **administration:**
528 Professor Darren L Smith PhD^{60, 61}.
529

- 530 **Project administration, metadata curation, samples, logistics, sequencing and analysis:**
531 Dr Martin D Curran PhD ¹⁴, and Dr Surendra Parmar PhD ¹⁴.
532
- 533 **Samples, logistics, metadata curation, project administration sequencing and analysis:**
534 Dr James G Shepherd MBChB MRCP ²¹.
535
- 536 **Sequencing, analysis, project administration, metadata curation and software and**
537 **analysis tools:**
538 Dr Matthew D Parker PhD ³⁸ and Dr Dinesh Aggarwal MRCP ^{1, 2, 3}.
539
- 540 **Leadership, supervision, funding acquisition, samples, logistics, and metadata curation:**
541 Dr Catherine Moore³³.
542
- 543 **Leadership, supervision, metadata curation, samples, logistics, sequencing and**
544 **analysis:**
545 Dr Derek J Fairley PhD ^{6, 88}, Professor Matthew W Loose PhD ⁵⁴, and Joanne Watkins MSc ³³.
546
- 547 **Metadata curation, sequencing, analysis, leadership, supervision and software and**
548 **analysis tools:**
549 Dr Matthew Bull PhD³³, and Dr Sam Nicholls PhD ¹⁵.
550
- 551 **Leadership, supervision, visualisation, sequencing, analysis and software and analysis**
552 **tools:**
553 Professor David M Aanensen PhD ^{1, 30}.
554
- 555 **Sequencing, analysis, samples, logistics, metadata curation, and visualisation:**
556 Dr Sharon Glaysher ⁷⁰.
557
- 558 **Metadata curation, sequencing, analysis, visualisation, software and analysis tools:**
559 Dr Matthew Bashton PhD ⁶⁰, and Dr Nicole Pacchiarini PhD ³³.
560
- 561 **Sequencing, analysis, visualisation, metadata curation, and software and analysis tools:**
562 Dr Anthony P Underwood PhD ^{1, 30}.
563
- 564 **Funding acquisition, leadership, supervision and project administration:**
565 Dr Thushan I de Silva PhD ³⁸, and Dr Dennis Wang PhD ³⁸.
566
- 567 **Project administration, samples, logistics, leadership and supervision:**
568 Dr Monique Andersson PhD ²⁸, Professor Anoop J Chauhan ⁷⁰, Dr Mariateresa de Cesare
569 PhD²⁶, Dr Catherine Ludden ^{1, 3}, and Dr Tabitha W Mahungu FRCPATH ⁹¹.
570
- 571 **Sequencing, analysis, project administration and metadata curation:**
572 Dr Rebecca Dewar PhD ²⁰, and Martin P McHugh MSc²⁰.
573

574 **Samples, logistics, metadata curation and project administration:**

575 Dr Natasha G Jesudason MBChB MRCP FRCPATH²¹, Dr Kathy K Li MBBCh FRCPATH²¹, Dr
576 Rajiv N Shah BMBS MRCP MSc²¹, and Dr Yusri Taha MD, PhD⁶⁶.

577

578 **Leadership, supervision, funding acquisition and metadata curation:**

579 Dr Kate E Templeton PhD²⁰.

580 **Leadership, supervision, funding acquisition, sequencing and analysis:**

581 Dr Simon Cottrell PhD³³, Dr Justin O'Grady PhD⁵¹, Professor Andrew Rambaut DPhil¹⁹, and
582 Professor Colin P Smith PhD⁹³.

583

584 **Leadership, supervision, metadata curation , sequencing and analysis:**

585 Professor Matthew T.G. Holden PhD⁸⁷, and Professor Emma C Thomson PhD/FRCP²¹.

586

587 **Leadership, supervision, samples, logistics and metadata curation:**

588 Dr Samuel Moses MD^{81, 82}.

589

590 **Sequencing, analysis, leadership, supervision, samples and logistics:**

591 Dr Meera Chand⁷, Dr Chrystala Constantinidou PhD⁷¹, Professor Alistair C Darby PhD⁴⁶,
592 Professor Julian A Hiscox PhD⁴⁶, Professor Steve Paterson PhD⁴⁶, and Dr Meera Unnikrishnan
593 PhD⁷¹.

594

595 **Sequencing, analysis, leadership and supervision and software and analysis tools:**

596 Dr Andrew J Page PhD⁵¹, and Dr Erik M Volz PhD⁹⁶.

597

598 **Samples, logistics, sequencing, analysis and metadata curation:**

599 Dr Charlotte J Houldcroft PhD⁸, Dr Aminu S Jahun PhD¹¹, Dr James P McKenna PhD⁸⁸, Dr
600 Luke W Meredith PhD¹¹, Dr Andrew Nelson PhD⁶¹, Sarojini Pandey MSc⁷², and Dr Gregory R
601 Young PhD⁶⁰.

602

603 **Sequencing, analysis, metadata curation, and software and analysis tools:**

604 Dr Anna Price PhD³⁴, Dr Sara Rey PhD³³, Dr Sunando Roy PhD⁴¹, Dr Ben Temperton Ph.D⁴⁹,
605 and Matthew Wyles³⁸.

606

607 **Sequencing, analysis, metadata curation and visualisation:**

608 Stefan Rooke MSc¹⁹, and Dr Sharif Shaaban PhD⁸⁷.

609

610 **Visualisation, sequencing, analysis and software and analysis tools:**

611 Dr Helen Adams PhD³⁵, Dr Yann Bourgeois Ph.D⁶⁹, Dr Katie F Loveson Ph.D⁶⁸, Áine O'Toole
612 MSc¹⁹, and Richard Stark MSc⁷¹.

613

614 **Project administration, leadership and supervision:**

615 Dr Ewan M Harrison PhD^{1, 3}, David Heyburn³³, and Professor Sharon J Peacock^{2, 3}

616

617 **Project administration and funding acquisition:**

618 Dr David Buck PhD²⁶, and Michaela John BSc Hons³⁶

619

620 **Sequencing, analysis and project administration:**

621 Dorota Jamrozy¹, and Dr Joshua Quick PhD¹⁵

622

623 **Samples, logistics, and project administration:**

624 Dr Rahul Batra MD⁷⁸, Katherine L Bellis BSc (Hons)^{1,3}, Beth Blane BSc³, Sophia T Girgis MSc³,
625 Dr Angie Green PhD²⁶, Anita Justice MSc²⁸, Dr Mark Kristiansen PhD⁴¹, and Dr Rachel J
626 Williams PhD⁴¹.

627

628 **Project administration, software and analysis tools:**

629 Radoslaw Poplawski BSc¹⁵.

630

631 **Project administration and visualisation:**

632 Dr Garry P Scarlett Ph.D⁶⁹.

633

634 **Leadership, supervision, and funding acquisition:**

635 Professor John A Todd PhD²⁶, Dr Christophe Fraser PhD²⁷, Professor Judith
636 Breuer MD^{40,41}, Professor Sergi Castellano PhD⁴¹, Dr Stephen L Michell PhD⁴⁹, Professor
637 Dimitris Gramatopoulos PhD, FRCPATH⁷³, and Dr Jonathan Edgeworth PhD, FRCPATH⁷⁸.

638

639 **Leadership, supervision and metadata curation:**

640 Dr Gemma L Kay PhD⁵¹.

641

642 **Leadership, supervision, sequencing and analysis:**

643 Dr Ana da Silva Filipe PhD²¹, Dr Aaron R Jeffries PhD⁴⁹, Dr Sascha Ott PhD⁷¹, Professor
644 Oliver Pybus²⁴, Professor David L Robertson PhD²¹, Dr David A Simpson PhD⁶, and Dr Chris
645 Williams MB BS³³.

646

647 **Samples, logistics, leadership and supervision:**

648 Dr Cressida Auckland FRCPATH⁵⁰, Dr John Boyes MBChB⁸³, Dr Samir Dervisevic FRCPATH⁵²,
649 Professor Sian Ellard FRCPATH^{49,50}, Dr Sonia Goncalves¹, Dr Emma J Meader FRCPATH⁵¹, Dr
650 Peter Muir PhD², Dr Husam Osman PhD⁹⁵, Reenesh Prakash MPH⁵², Dr Venkat Sivaprakasam
651 PhD¹⁸, and Dr Ian B Vipond PhD².

652

653 **Leadership, supervision and visualisation**

654 Dr Jane AH Masoli MBChB^{49,50}.

655

656 **Sequencing, analysis and metadata curation**

657 Dr Nabil-Fareed Alikhan PhD⁵¹, Matthew Carlile BSc⁵⁴, Dr Noel Craine DPhil³³, Dr Sam T
658 Haldenby PhD⁴⁶, Dr Nadine Holmes PhD⁵⁴, Professor Ronan A Lyons MD³⁷, Dr Christopher
659 Moore PhD⁵⁴, Malorie Perry MSc³³, Dr Ben Warne MRCP⁸⁰, and Dr Thomas Williams
660 MD¹⁹.

661 **Samples, logistics and metadata curation:**

662 Dr Lisa Berry PhD ⁷², Dr Andrew Bosworth PhD ⁹⁵, Dr Julianne Rose Brown PhD⁴⁰, Sharon
 663 Campbell MSc⁶⁷, Dr Anna Casey PhD ¹⁷, Dr Gemma Clark PhD ⁵⁶, Jennifer Collins BSc ⁶⁶,
 664 Dr Alison Cox PhD ^{43, 44}, Thomas Davis MSc ⁸⁴, Gary Eltringham BSc ⁶⁶, Dr Cariad Evans ^{38, 39},
 665 Dr Clive Graham MD ⁶⁴, Dr Fenella Halstead PhD ¹⁸, Dr Kathryn Ann Harris PhD ⁴⁰, Dr
 666 Christopher Holmes PhD ⁵⁸, Stephanie Hutchings ², Professor Miren Iturriza-Gomara PhD ⁴⁶,
 667 Dr Kate Johnson ^{38, 39}, Katie Jones MSc ⁷², Dr Alexander J Keeley MRCP ³⁸, Dr Bridget A
 668 Knight PhD ^{49, 50}, Cherian Koshy MSc, CSci, FIBMS ⁹⁰, Steven Liggett ⁶³, Hannah Lowe MSc
 669 ⁸¹, Dr Anita O Lucaci PhD ⁴⁶, Dr Jessica Lynch PhD MBChB ^{25, 29}, Dr Patrick C McClure PhD ⁵⁵
 670, Dr Nathan Moore MBChB ³¹, Matilde Mori BSc ^{25, 29, 32}, Dr David G Partridge FRCP, FRCPath
 671 ^{38, 39}, Pinglawathee Madona ^{43, 44}, Hannah M Pymont MSc ², Dr Paul Anthony Randell MBBCh
 672 ^{43, 44}, Dr Mohammad Raza ^{38, 39}, Felicity Ryan MSc ⁸¹, Dr Robert Shaw FRCPath ²⁸, Dr Tim J
 673 Sloan PhD ⁵⁷, and Emma Swindells BSc ⁶⁵.

674

675 **Sequencing, analysis, Samples and logistics:**

676 Alexander Adams BSc ³³, Dr Hibo Asad PhD ³³, Alec Birchley MSc ³³, Tony Thomas
 677 Brooks BSc (Hons) ⁴¹, Dr Giselda Bucca PhD ⁹³, Ethan Butcher ⁷⁰, Dr Sarah L Caddy PhD ¹³, Dr
 678 Laura G Caller PhD ^{2, 3, 12}, Yasmin Chaudhry BSc ¹¹, Jason Coombes BSc (HONS) ³³, Michelle
 679 Cronin ³³, Patricia L Dyal MPhil ⁴¹, Johnathan M Evans MSc ³³, Laia Fina ³³, Bree Gatica-Wilcox
 680 MPhil ³³, Dr Iliana Georgana PhD ¹¹, Lauren Gilbert A-Levels ³³, Lee Graham BSc ³³, Danielle C
 681 Groves BA ³⁸, Grant Hall BSc ¹¹, Ember Hilvers MPH ³³, Dr Myra Hosmillo PhD ¹¹,
 682 Hannah Jones ³³, Sophie Jones MSc ³³, Fahad A Khokhar BSc ¹³, Sara Kumziene-
 683 Summerhayes MSc ³³, George MacIntyre-Cockett BSc ²⁶, Dr Rocio T Martinez Nunez PhD ⁹⁴,
 684 Dr Caoimhe McKerr PhD ³³, Dr Claire McMurray PhD ¹⁵, Dr Richard Myers ⁷, Yasmin Nicole
 685 Panchbhaya BSc ⁴¹, Malte L Pinckert MPhil ¹¹, Amy Plimmer ³³, Dr Joanne Stockton PhD ¹⁵,
 686 Sarah Taylor ³³, Dr Alicia Thornton ⁷, Amy Trebes MSc ²⁶, Alexander J Trotter MRes ⁵¹
 687, Helena Jane Tutill BSc ⁴¹, Charlotte A Williams BSc ⁴¹, Anna Yakovleva BSc ¹¹ and Dr Wen C
 688 Yew PhD ⁶².

689

690 **Sequencing, analysis and software and analysis tools:**

691 Dr Mohammad T Alam PhD ⁷¹, Dr Laura Baxter PhD ⁷¹, Olivia Boyd MSc ⁹⁶, Dr Fabricia
 692 F. Nascimento PhD ⁹⁶, Timothy M Freeman MPhil ³⁸, Lily Geidelberg MSc ⁹⁶, Dr Joseph Hughes
 693 PhD ²¹, David Jorgensen MSc ⁹⁶, Dr Benjamin B Lindsey MRCP ³⁸, Dr Richard J Orton PhD ²¹,
 694 Dr Manon Ragonnet-Cronin PhD ⁹⁶, Joel Southgate MSc ^{33, 34}, and Dr Sreenu Vattipally PhD ²¹.

695

696 **Samples, logistics and software and analysis tools:**

697 Dr Igor Starinskij MSc MRCP ²³.

698

699 **Visualisation and software and analysis tools:**

700 Dr Joshua B Singer PhD ²¹, Dr Khalil Abudahab PhD ^{1, 30}, Leonardo de Oliveira Martins PhD⁵¹,
 701 Dr Thanh Le-Viet PhD ⁵¹, Mirko Menegazzo ³⁰, Ben EW Taylor Meng ^{1, 30}, and Dr Corin A
 702 Yeats PhD ³⁰.

703

704 **Project Administration:**

705 Sophie Palmer³, Carol M Churcher³, Dr Alisha Davies³³, Elen De Lacy MSc³³, Fatima
 706 Downing³³, Sue Edwards³³, Dr Nikki Smith PhD³⁸, Dr Francesc Coll PhD⁹⁷, Dr
 707 Nazreen F Hadjirin PhD³ and Dr Frances Bolt PhD^{44, 45}.

708

709 **Leadership and supervision:**

710 Dr. Alex Alderton¹, Dr Matt Berriman¹, Ian G Charles⁵¹, Dr Nicholas Cortes MBChB³¹, Dr
 711 Tanya Curran PhD⁸⁸, Prof John Danesh¹, Dr Sahar Eldirdiri MBBS, MSc FRCPATH⁸⁴, Dr
 712 Ngozi Elumogo FRCPATH⁵², Prof Andrew Hattersley FRS^{49, 50}, Professor Alison Holmes MD^{44,}
 713 ⁴⁵, Dr Robin Howe³³, Dr Rachel Jones³³, Anita Kenyon MSc⁸⁴, Prof Robert A Kingsley PhD⁵¹,
 714 Professor Dominic Kwiatkowski^{1, 9}, Dr Cordelia Langford¹, Dr Jenifer Mason MBBS⁴⁸, Dr Alison
 715 E Mather PhD⁵¹, Lizzie Meadows MA⁵¹, Dr Sian Morgan FRCPATH³⁶, Dr James Price PhD^{44,}
 716 ⁴⁵, Trevor I Robinson MSc⁴⁸, Dr Giri Shankar³³, John Wain⁵¹, and Dr Mark A Webber PhD⁵¹

717

718

719 **Metadata curation:**

720 Dr Declan T Bradley PhD^{5, 6}, Dr Michael R Chapman PhD^{1, 3, 4}, Dr Derrick Croke²⁸, Dr David
 721 Eyre PhD²⁸, Professor Martyn Guest PhD³⁴, Huw Gulliver³⁴, Dr Sarah Hoosdally²⁸, Dr
 722 Christine Kitchen PhD³⁴, Dr Ian Merrick PhD³⁴, Siddharth Mookerjee MPH^{44, 45}, Robert Munn
 723 BSc³⁴, Professor Timothy Peto PhD²⁸, Will Potter⁵², Dr Dheeraj K Sethi MBBS⁵²,
 724 Wendy Smith⁵⁶, Dr Luke B Snell MB BS^{75, 94}, Dr Rachael Stanley PhD⁵², Claire Stuart⁵² and
 725 Dr Elizabeth Wastenge MD²⁰.

726

727 **Sequencing and analysis:**

728 Dr Erwan Acheson PhD⁶, Safiah Afifi BSc³⁶, Dr Elias Allara MD PhD^{2, 3}, Dr Roberto
 729 Amato¹, Dr Adrienn Angyal PhD³⁸, Dr Elihu Aranday-Cortes PhD/DVM²¹, Cristina Ariani¹,
 730 Jordan Ashworth¹⁹, Dr Stephen Attwood²⁴, Alp Aydin MSci⁵¹, David J Baker BEng⁵¹, Dr
 731 Carlos E Balcazar PhD¹⁹, Angela Beckett MSc⁶⁸, Robert Beer BSc³⁶, Dr Gilberto
 732 Betancor PhD⁷⁶, Emma Betteridge¹, Dr David Bibby⁷, Dr Daniel Bradshaw⁷,
 733 Catherine Bresner BSc(Hons)³⁴, Dr Hannah E Bridgewater PhD⁷¹, Alice Broos BSc (Hons)²¹,
 734 Dr Rebecca Brown PhD³⁸, Dr Paul E Brown PhD⁷¹, Dr Kirstyn Brunker PhD²², Dr Stephen N
 735 Carmichael PhD²¹, Jeffrey K. J. Cheng MSc⁷¹, Dr Rachel Colquhoun DPhil¹⁹, Dr Gavin
 736 Dabrera⁷, Dr Johnny Debebe PhD⁵⁴, Eleanor Drury¹, Dr Louis du Plessis²⁴, Richard Eccles
 737 MSc⁴⁶, Dr Nicholas Ellaby⁷, Audrey Farbos MSc⁴⁹, Ben Farr¹, Dr Jacqueline Findlay PhD⁴¹,
 738 Chloe L Fisher MSc⁷⁴, Leysa Marie Forrest MSc⁴¹, Dr Sarah Francois²⁴, Lucy R. Frost BSc⁷¹,
 739 William Fuller BSc³⁴, Dr Eileen Gallagher⁷, Dr Michael D Gallagher PhD¹⁹, Matthew Gemmell
 740 MSc⁴⁶, Dr Rachel AJ Gilroy PhD⁵¹, Scott Goodwin¹, Dr Luke R Green PhD³⁸, Dr Richard
 741 Gregory PhD⁴⁶, Dr Natalie Groves⁷, Dr James W Harrison PhD⁴⁹, Hassan Hartman⁷, Dr
 742 Andrew R Hesketh PhD⁹³, Verity Hill¹⁹, Dr Jonathan Hubb⁷, Dr Margaret Hughes PhD⁴⁶, Dr
 743 David K Jackson¹, Dr Ben Jackson PhD¹⁹, Dr Keith James¹, Natasha Johnson BSc (Hons)²¹
 744 , Ian Johnston¹, Jon-Paul Keatley¹, Dr Moritz Kraemer²⁴, Dr Angie Lackenby⁷, Dr Mara
 745 Lawniczak¹, Dr David Lee⁷, Rich Livett¹, Stephanie Lo¹, Daniel Mair BSc (Hons)²¹, Joshua
 746 Maksimovic FD sport science³⁶, Nikos Manesis⁷, Dr Robin Manley Ph.D⁴⁹, Dr Carmen Manso⁷
 747 , Dr Angela Marchbank BSc³⁴, Dr Inigo Martincorena¹, Dr Tamyo Mbisa⁷, Kathryn McCluggage
 748 MSC³⁶, Dr JT McCrone PhD¹⁹, Shahjahan Miah⁷, Michelle L Michelsen BSc⁴⁹, Dr Mari

749 Morgan PhD³³, Dr Gaia Nebbia PhD, FRCPath⁷⁸, Charlotte Nelson MSc⁴⁶, Jenna Nichols BSc
 750 (Hons)²¹, Dr Paola Niola PhD⁴¹, Dr Kyriaki Nomikou PhD²¹, Steve Palmer¹, Dr. Naomi Park¹,
 751 Dr Yasmin A Parr PhD²¹, Dr Paul J Parsons PhD³⁸, Vineet Patel⁷, Dr. Minal Patel¹, Clare
 752 Pearson MSc^{2,1}, Dr Steven Platt⁷, Christoph Puethe¹, Dr. Mike Quail¹, Dr JaynaRaghwani²⁴,
 753 Dr Lucille Rainbow PhD⁴⁶, Shavanthi Rajatileka¹, Dr Mary Ramsay⁷, Dr Paola C Resende
 754 Silva PhD^{41,42}, Steven Rudder⁵¹, Dr Chris Ruis³, Dr Christine M Sambles PhD⁴⁹, Dr Fei Sang
 755 PhD⁵⁴, Dr Ulf Schaefer⁷, Dr Emily Scher PhD¹⁹, Dr. Carol Scott¹, Lesley Shirley¹, Adrian W
 756 Signell BSc⁷⁶, John Sillitoe¹, Christen Smith¹, Dr Katherine L Smollett PhD²¹, Karla Spellman
 757 FD³⁶, Thomas D Stanton BSc¹⁹, Dr David J Studholme PhD⁴⁹, Ms Grace Taylor-Joyce BSc⁷¹
 758 ,Dr Ana P Tedim PhD⁵¹, Dr Thomas Thompson PhD⁶, Dr Nicholas M Thomson PhD⁵¹, Scott
 759 Thurston¹, Lily Tong PhD²¹, Gerry Tonkin-Hill¹, Rachel M Tucker MSc³⁸, Dr Edith E Vamos
 760 PhD⁴, Dr Tetyana Vasylyeva²⁴, Joanna Warwick-Dugdale BSc⁴⁹, Danni Weldon¹, Dr Mark
 761 Whitehead PhD⁴⁶, Dr David Williams⁷, Dr Kathleen A Williamson PhD¹⁹, Harry D Wilson BSc
 762 ⁷⁶, Trudy Workman HNC³⁴, Dr Muhammad Yasir PhD⁵¹, Dr Xiaoyu Yu PhD¹⁹, and Dr Alex
 763 Zarebski²⁴.

764

765 **Samples and logistics:**

766 Dr Evelien M Adriaenssens PhD⁵¹, Dr Shazaad S Y Ahmad MSc^{2,47}, Adela Alcolea-Medina
 767 MPharm^{59,77}, Dr John Allan PhD⁶⁰, Dr Patawee Asamaphan PhD²¹, Laura Atkinson MSc⁴⁰,
 768 Paul Baker MD⁶³, Professor Jonathan Ball PhD⁵⁵, Dr Edward Barton MD⁶⁴, Dr. Mathew A
 769 Beale¹, Dr. Charlotte Beaver¹, Dr Andrew Beggs PhD¹⁶, Dr Andrew Bell PhD⁵¹, Duncan J
 770 Berger¹, Dr Louise Berry⁵⁶, Claire M Bewshea MSc⁴⁹, Kelly Bicknell⁷⁰, Paul Bird⁵⁸, Dr Chloe
 771 Bishop⁷, Dr Tim Boswell⁵⁶, Cassie Breen BSc⁴⁸, Dr Sarah K Buddenborg¹, Dr Shirelle Burton-
 772 Fanning MD⁶⁶, Dr Vicki Chalker⁷, Dr Joseph G Chappell PhD⁵⁵, Themoula Charalampous MSc
 773 ^{78,94}, Claire Cormie³, Dr Nick Cortes PhD^{29,25}, Dr Lindsay J Coupland PhD⁵², Angela Cowell
 774 MSc⁴⁸, Dr Rose K Davidson PhD⁵³, Joana Dias MSc³, Dr Maria Diaz PhD⁵¹, Thomas Dibling¹,
 775 Matthew J Dorman¹, Dr Nichola Duckworth⁵⁷, Scott Elliott⁷⁰, Sarah Essex⁶³, Karlie Fallon⁵⁸,
 776 Theresa Feltwell⁸, Dr Vicki M Fleming PhD⁵⁶, Sally Forrest BSc³, Luke Foulser¹, Maria V
 777 Garcia-Casado¹, Dr Artemis Gavriil PhD⁴¹, Dr Ryan P George PhD⁴⁷, Laura Gifford MSc³³,
 778 Harmeet K Gill PhD³, Jane Greenaway MSc⁶⁵, Luke Griffith BSc⁵³, Ana Victoria Gutierrez⁵¹, Dr
 779 Antony D Hale MBBS⁸⁵, Dr Tanzina Haque FRCPath, PhD⁹¹, Katherine L Harper MBiol⁸⁵, Dr Ian
 780 Harrison⁷, Dr Judith Heaney PhD⁸⁹, Thomas Helmer⁵⁸, Ellen E Higginson PhD³, Richard
 781 Hopes², Dr Hannah C Howson-Wells PhD⁵⁶, Dr Adam D Hunter¹, Robert Impey⁷⁰, Dr Dianne
 782 Irish-Tavares FRCPath⁹¹, David A Jackson¹, Kathryn A Jackson MSc⁴⁶, Dr Amelia Joseph⁵⁶,
 783 Leanne Kane¹, Sally Kay¹, Leanne M Kermack MSc³, Manjinder Khakh⁵⁶, Dr Stephen P Kidd
 784 PhD^{29,25,31}, Dr Anastasia Kolyva PhD⁵¹, Jack CD Lee BSc⁴⁰, Laura Letchford¹, Nick Levene
 785 MSc⁷⁹, Dr LisaJ Levett PhD⁸⁹, Dr Michelle M Lister PhD⁵⁶, Allyson Lloyd⁷⁰, Dr Joshua Loh
 786 PhD⁶⁰, Dr Louissa R Macfarlane-Smith PhD⁸⁵, Dr Nicholas W Machin MSc^{2,47}, Mailis Maes
 787 M.phil³, Dr Samantha McGuigan¹, Liz McMinn¹, Dr Lamia Mestek-Boukhibar D.Phil⁴¹, Dr
 788 Zoltan Molnar PhD⁶, Lynn Monaghan⁷⁹, Dr Catrin Moore²⁷, Plamena Naydenova BSc³,
 789 Alexandra S Neaverson¹, Dr. Rachel Nelson PhD¹, Marc O Niebel MSc²¹, Elaine O'Toole BSc
 790 ⁴⁸, Debra Padgett BSc⁶⁴, Gaurang Patel¹, Dr Brendan Al Payne MD⁶⁶, Liam Prestwood¹, Dr
 791 Veena Raviprakash MD⁶⁷, Nicola Reynolds PhD⁸⁶, Dr Alex Richter PhD¹⁶, Dr Esther Robinson
 792 PhD⁹⁵, Dr Hazel A Rogers¹, Dr Aileen Rowan PhD⁹⁶, Garren Scott BSc⁶⁴, Dr Divya Shah

793 PhD⁴⁰, Nicola Sheriff BSc⁶⁷, Dr Graciela Sluga MD - MSc⁹², Emily Souster¹, Dr. Michael
 794 Spencer-Chapman¹, Sushmita Sridhar BSc^{1,3}, Tracey Swingler⁵³, Dr Julian Tang⁵⁸, Professor
 795 Graham P Taylor DSc⁹⁶, Dr Theocharis Tsoleridis PhD⁵⁵, Dr Lance Turtle PhD MRCP⁴⁶, Dr
 796 Sarah Walsh⁵⁷, Dr Michelle Wantoch PhD⁸⁶, Joanne Watts BSc⁴⁸, Dr Sheila Waugh MD⁶⁶,
 797 Sam Weeks⁴¹, Dr Rebecca Williams BMBS³¹, Dr Iona Willingham⁵⁶, Dr Emma L Wise PhD^{25,}
 798 ^{29,31}, Victoria Wright BSc⁵⁴, Dr Sarah Wyllie⁷⁰, and Jamie Young BSc³.

799

800 **Software and analysis tools**

801 Amy Gaskin MSc³³, Dr Will Rowe PhD¹⁵, and Dr Igor Siveroni PhD⁹⁶.

802

803 **Visualisation:**

804 Dr Robert Johnson PhD⁹⁶.

805

806 **1** Wellcome Sanger Institute, **2** Public Health England, **3** University of Cambridge, **4** Health Data
 807 Research UK, Cambridge, **5** Public Health Agency, Northern Ireland, **6** Queen's University Belfast **7**
 808 Public Health England Colindale, **8** Department of Medicine, University of Cambridge, **9** University of
 809 Oxford, **10** Departments of Infectious Diseases and Microbiology, Cambridge University Hospitals NHS
 810 Foundation Trust; Cambridge, UK, **11** Division of Virology, Department of Pathology, University of
 811 Cambridge, **12** The Francis Crick Institute, **13** Cambridge Institute for Therapeutic Immunology and
 812 Infectious Disease, Department of Medicine, **14** Public Health England, Clinical Microbiology and Public
 813 Health Laboratory, Cambridge, UK, **15** Institute of Microbiology and Infection, University of Birmingham,
 814 **16** University of Birmingham, **17** Queen Elizabeth Hospital, **18** Heartlands Hospital, **19** University of
 815 Edinburgh, **20** NHS Lothian, **21** MRC-University of Glasgow Centre for Virus Research, **22** Institute of
 816 Biodiversity, Animal Health & Comparative Medicine, University of Glasgow, **23** West of Scotland
 817 Specialist Virology Centre, **24** Dept Zoology, University of Oxford, **25** University of Surrey, **26** Wellcome
 818 Centre for Human Genetics, Nuffield Department of Medicine, University of Oxford, **27** Big Data Institute,
 819 Nuffield Department of Medicine, University of Oxford, **28** Oxford University Hospitals NHS Foundation
 820 Trust, **29** Basingstoke Hospital, **30** Centre for Genomic Pathogen Surveillance, University of Oxford, **31**
 821 Hampshire Hospitals NHS Foundation Trust, **32** University of Southampton, **33** Public Health Wales NHS
 822 Trust, **34** Cardiff University, **35** Betsi Cadwaladr University Health Board, **36** Cardiff and Vale University
 823 Health Board, **37** Swansea University, **38** University of Sheffield, **39** Sheffield Teaching Hospitals, **40**
 824 Great Ormond Street NHS Foundation Trust, **41** University College London, **42** Oswaldo Cruz Institute,
 825 Rio de Janeiro **43** North West London Pathology, **44** Imperial College Healthcare NHS Trust, **45** NIHR
 826 Health Protection Research Unit in HCAI and AMR, Imperial College London, **46** University of Liverpool,
 827 **47** Manchester University NHS Foundation Trust, **48** Liverpool Clinical Laboratories, **49** University of
 828 Exeter, **50** Royal Devon and Exeter NHS Foundation Trust, **51** Quadram Institute Bioscience, University
 829 of East Anglia, **52** Norfolk and Norwich University Hospital, **53** University of East Anglia, **54** Deep Seq,
 830 School of Life Sciences, Queens Medical Centre, University of Nottingham, **55** Virology, School of Life
 831 Sciences, Queens Medical Centre, University of Nottingham, **56** Clinical Microbiology Department,
 832 Queens Medical Centre, **57** PathLinks, Northern Lincolnshire & Goole NHS Foundation Trust, **58** Clinical
 833 Microbiology, University Hospitals of Leicester NHS Trust, **59** Viapath, **60** Hub for Biotechnology in the
 834 Built Environment, Northumbria University, **61** NU-OMICS Northumbria University, **62** Northumbria
 835 University, **63** South Tees Hospitals NHS Foundation Trust, **64** North Cumbria Integrated Care NHS
 836 Foundation Trust, **65** North Tees and Hartlepool NHS Foundation Trust, **66** Newcastle Hospitals NHS
 837 Foundation Trust, **67** County Durham and Darlington NHS Foundation Trust, **68** Centre for Enzyme
 838 Innovation, University of Portsmouth, **69** School of Biological Sciences, University of Portsmouth, **70**
 839 Portsmouth Hospitals NHS Trust, **71** University of Warwick, **72** University Hospitals Coventry and

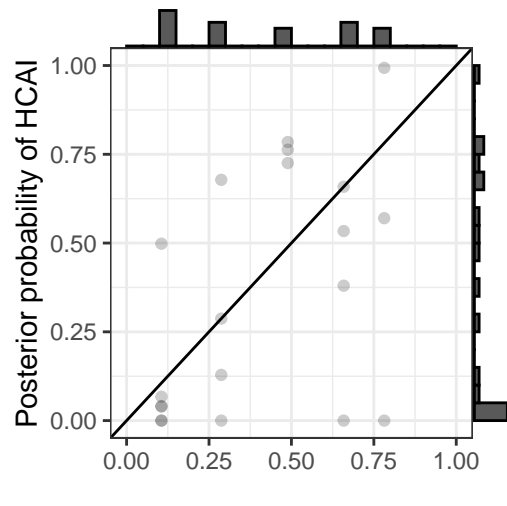
840 Warwickshire, **73** Warwick Medical School and Institute of Precision Diagnostics, Pathology, UHCW NHS
841 Trust, **74** Genomics Innovation Unit, Guy's and St. Thomas' NHS Foundation Trust, **75** Centre for Clinical
842 Infection & Diagnostics Research, St. Thomas' Hospital and Kings College London, **76** Department of
843 Infectious Diseases, King's College London, **77** Guy's and St. Thomas' Hospitals NHS Foundation Trust,
844 **78** Centre for Clinical Infection and Diagnostics Research, Department of Infectious Diseases, Guy's and
845 St Thomas' NHS Foundation Trust, **79** Princess Alexandra Hospital Microbiology Dept. , **80** Cambridge
846 University Hospitals NHS Foundation Trust, **81** East Kent Hospitals University NHS Foundation Trust, **82**
847 University of Kent, **83** Gloucestershire Hospitals NHS Foundation Trust, **84** Department of Microbiology,
848 Kettering General Hospital, **85** National Infection Service, PHE and Leeds Teaching Hospitals Trust, **86**
849 Cambridge Stem Cell Institute, University of Cambridge, **87** Public Health Scotland, 88 Belfast Health &
850 Social Care Trust, **89** Health Services Laboratories, **90** Barking, Havering and Redbridge University
851 Hospitals NHS Trust, **91** Royal Free NHS Trust, **92** Maidstone and Tunbridge Wells NHS Trust, **93**
852 University of Brighton, **94** Kings College London, **95** PHE Heartlands, **96** Imperial College London, **97**
853 Department of Infection Biology, London School of Hygiene and Tropical Medicine.

854

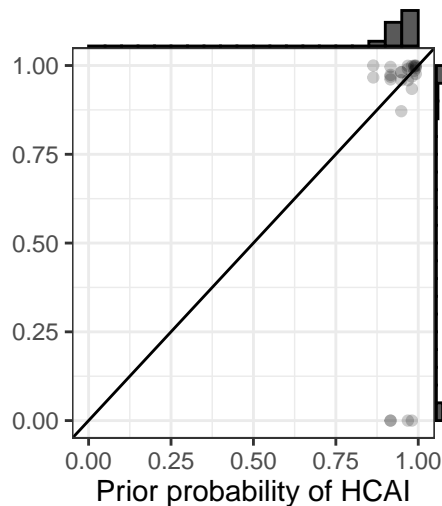
855

856

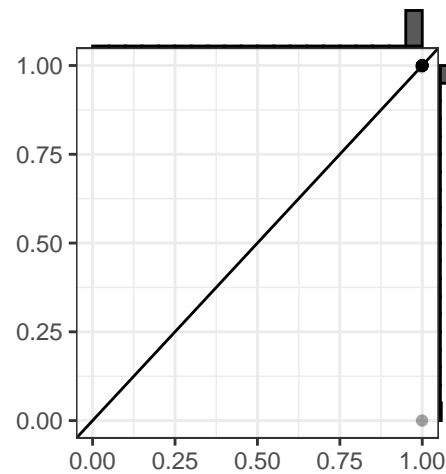
(a) Indeterminate HCAI



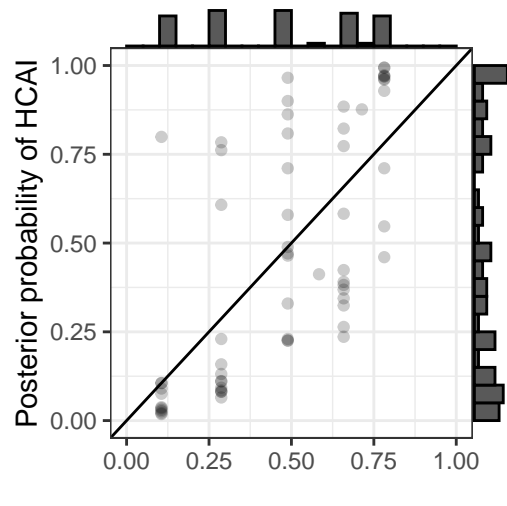
Probable HCAI



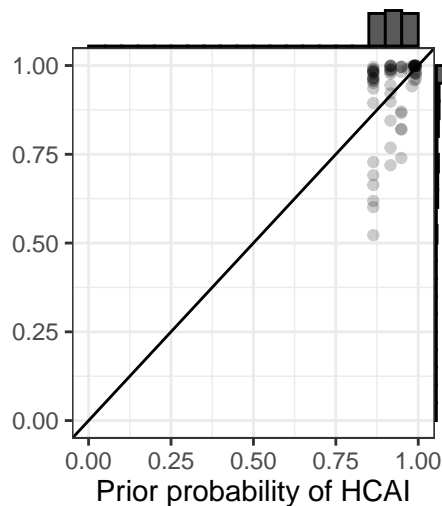
Definite HCAI



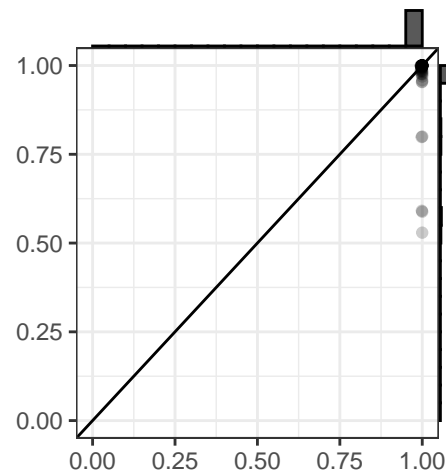
(b) Indeterminate HCAI

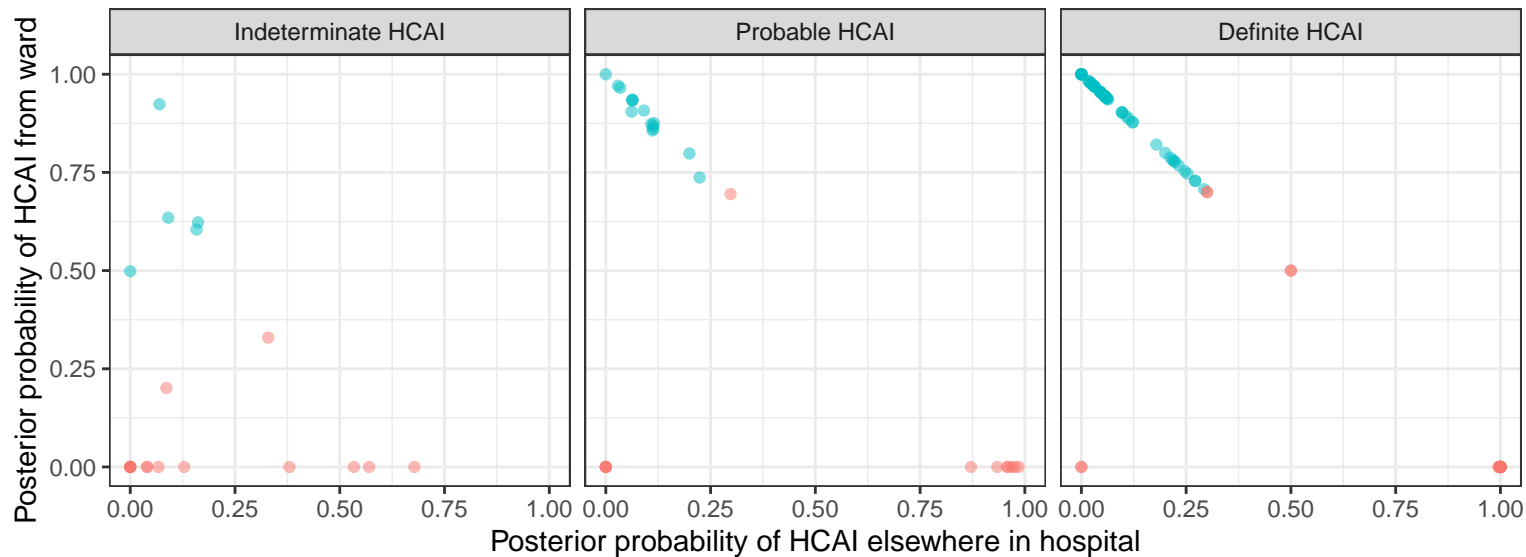
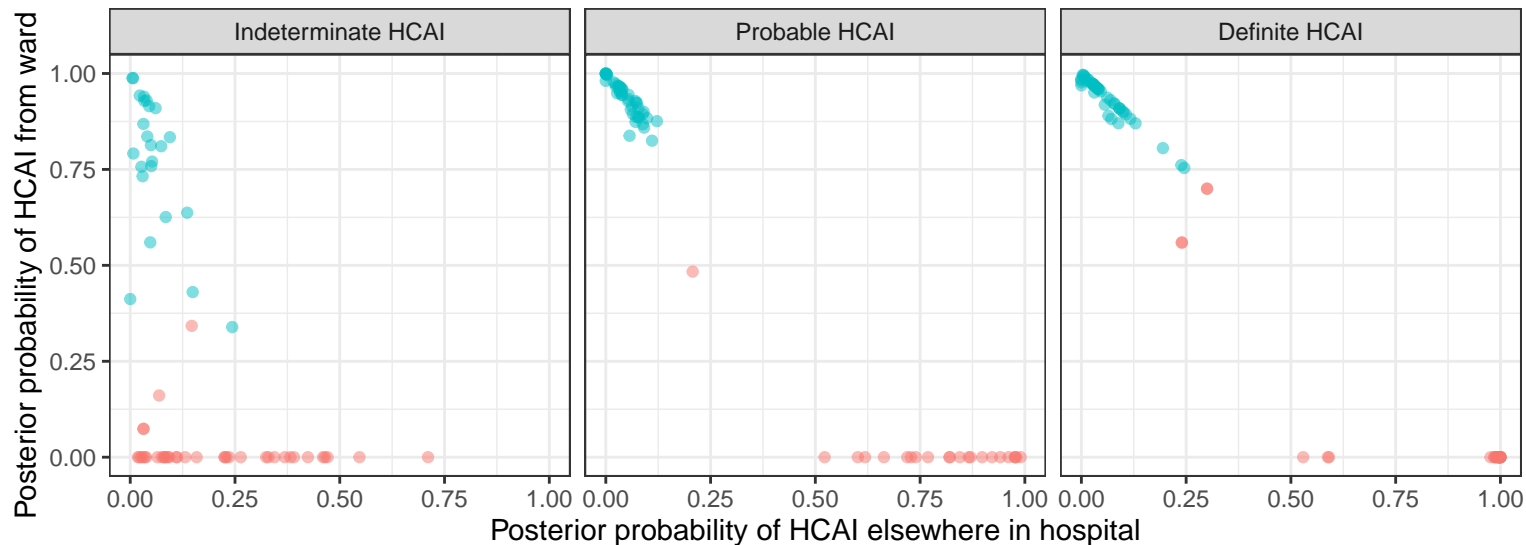


Probable HCAI



Definite HCAI



(a)**(b)**

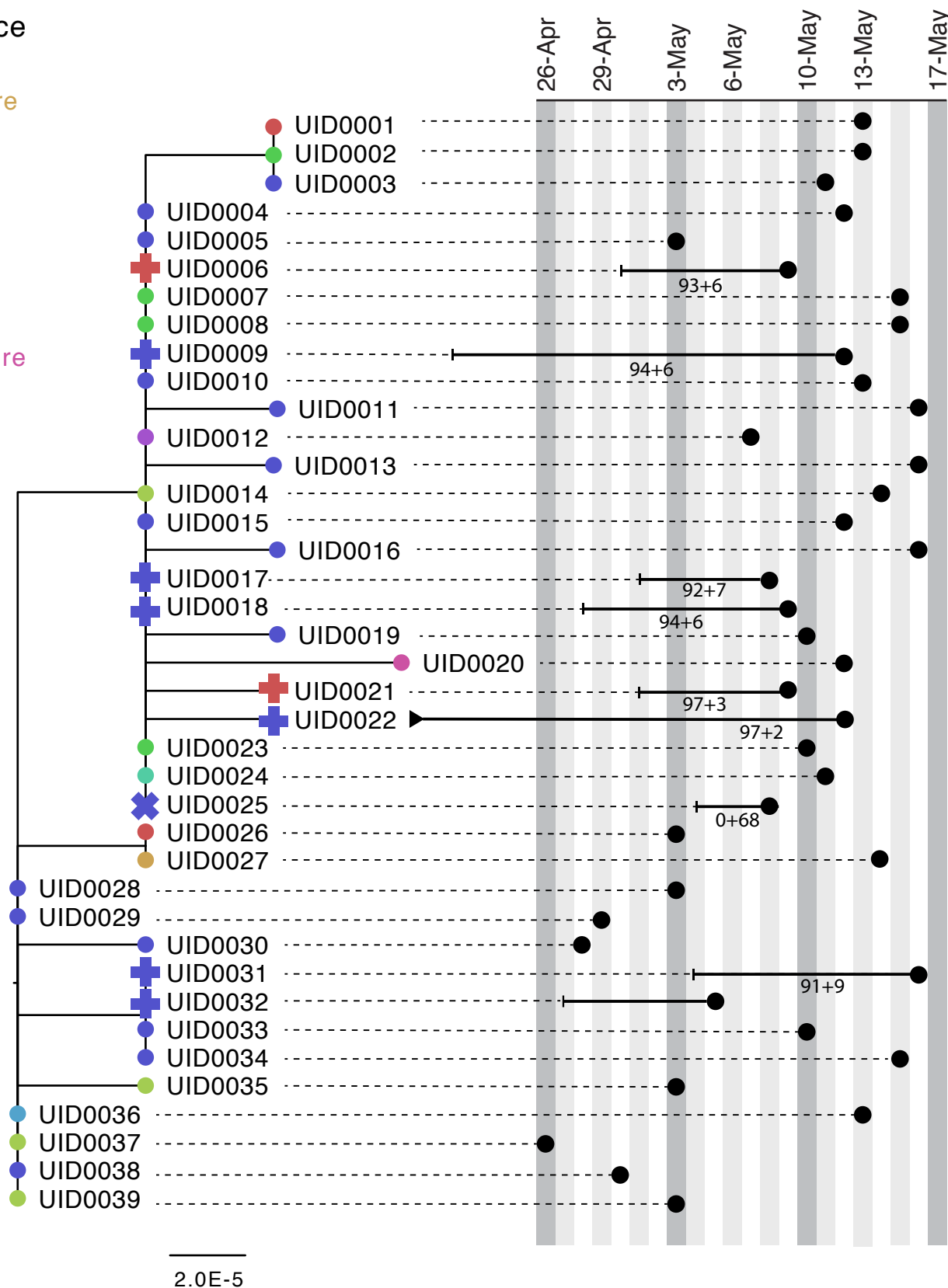
Sequence matches on ward ● None ● One or more

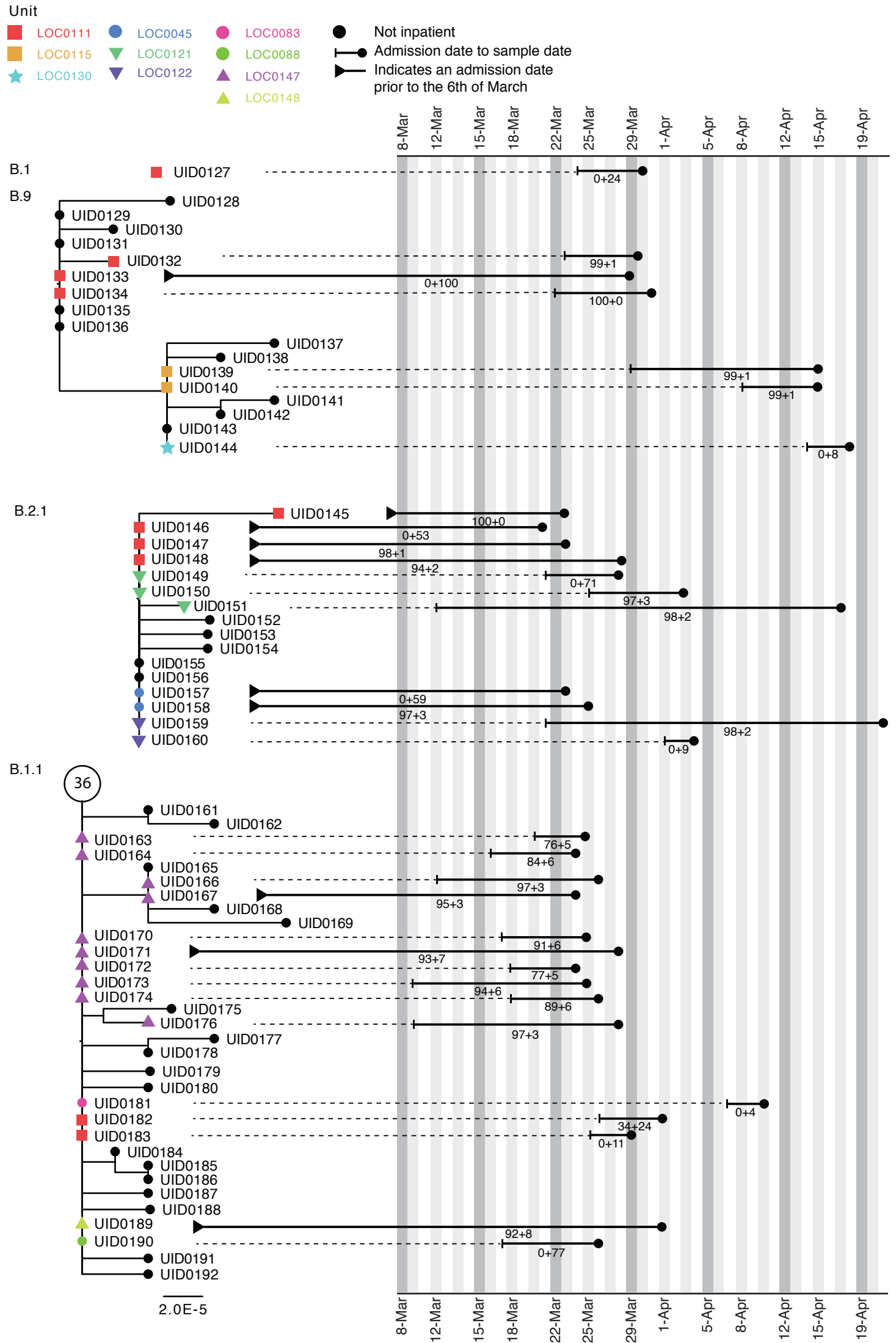
Location of residence

- Argyll and Bute
- East Dunbartonshire
- East Renfrewshire
- Glasgow City
- Inverclyde
- North Ayrshire
- Renfrewshire
- South Ayrshire
- West Dunbartonshire

Hospital 5 units

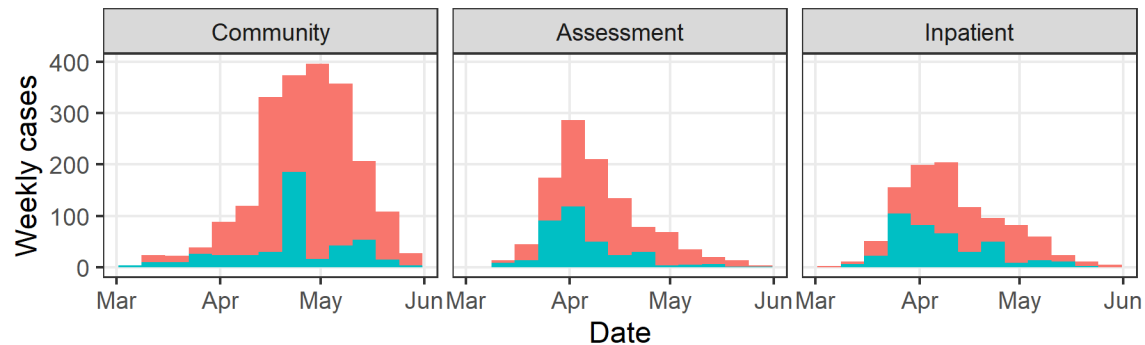
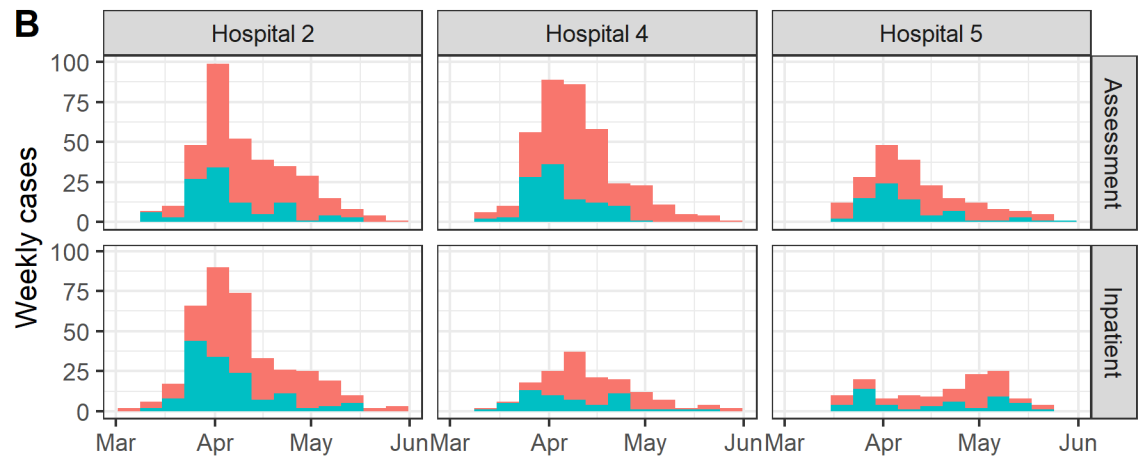
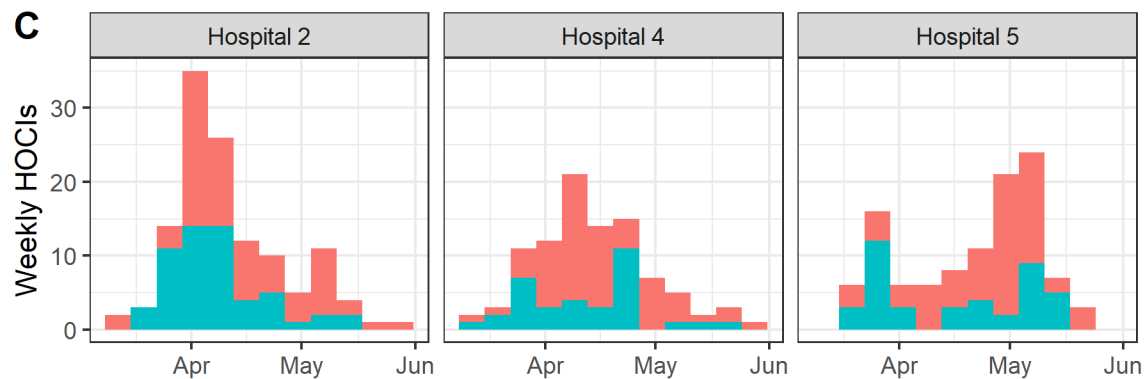
- ⊕ Unit 93
- ⊗ Unit 92



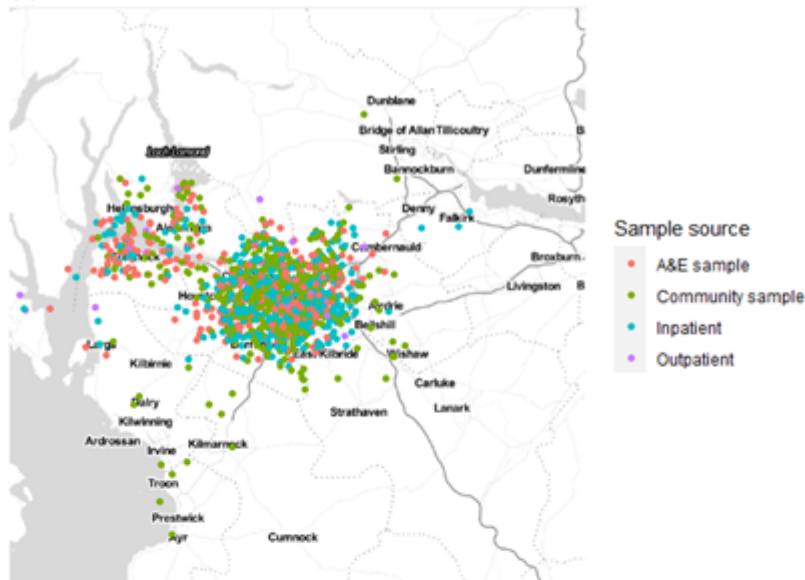


A

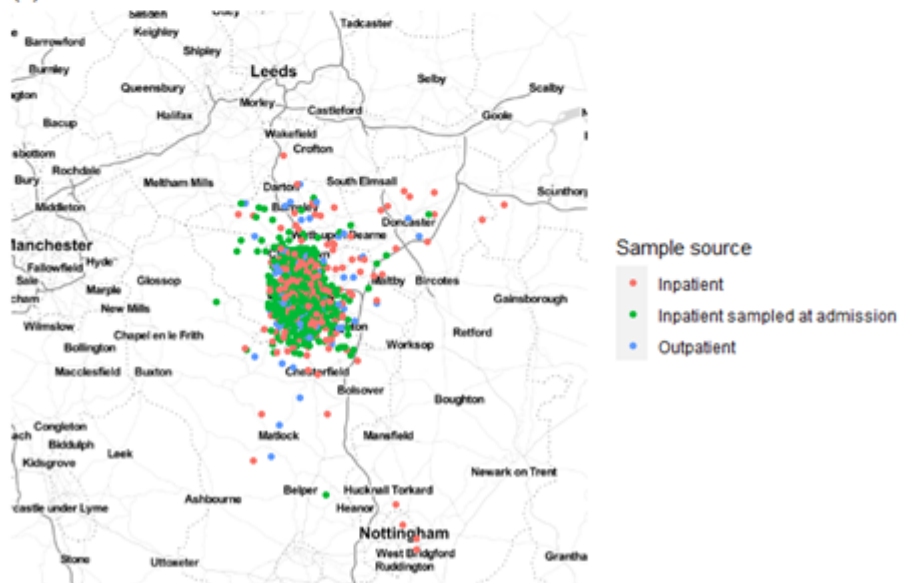
Sequenced Not sequenced

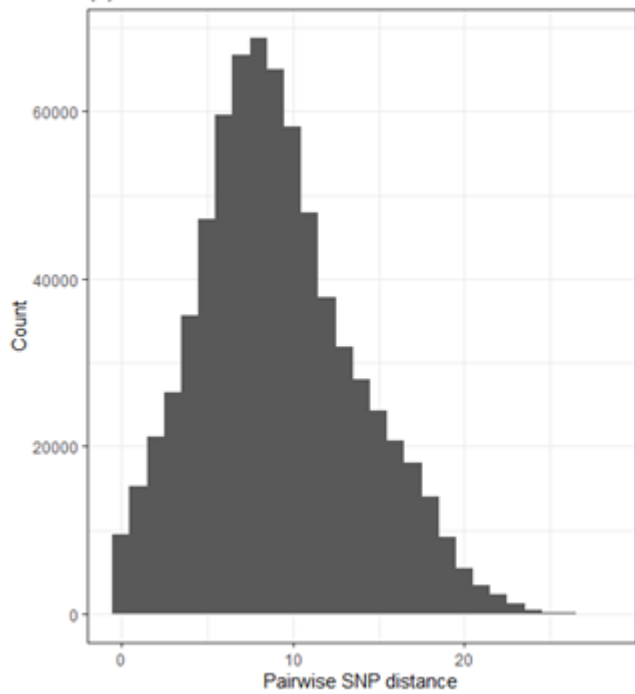
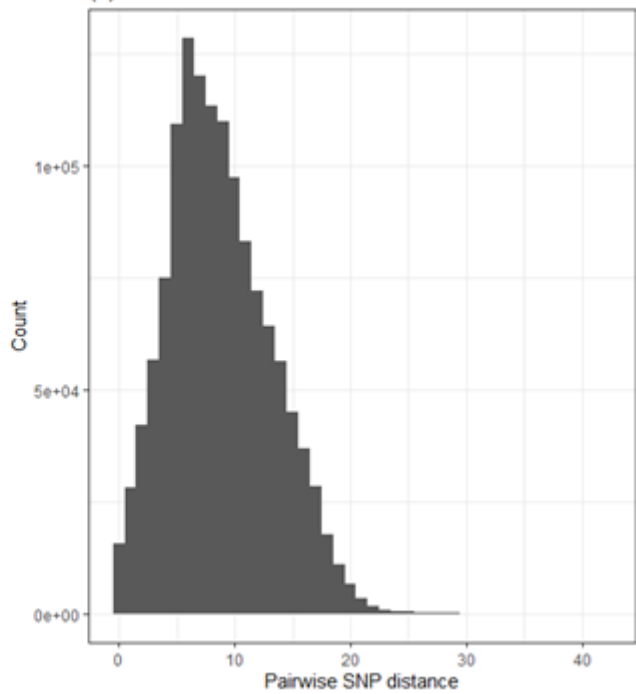
**B****C**

(a)

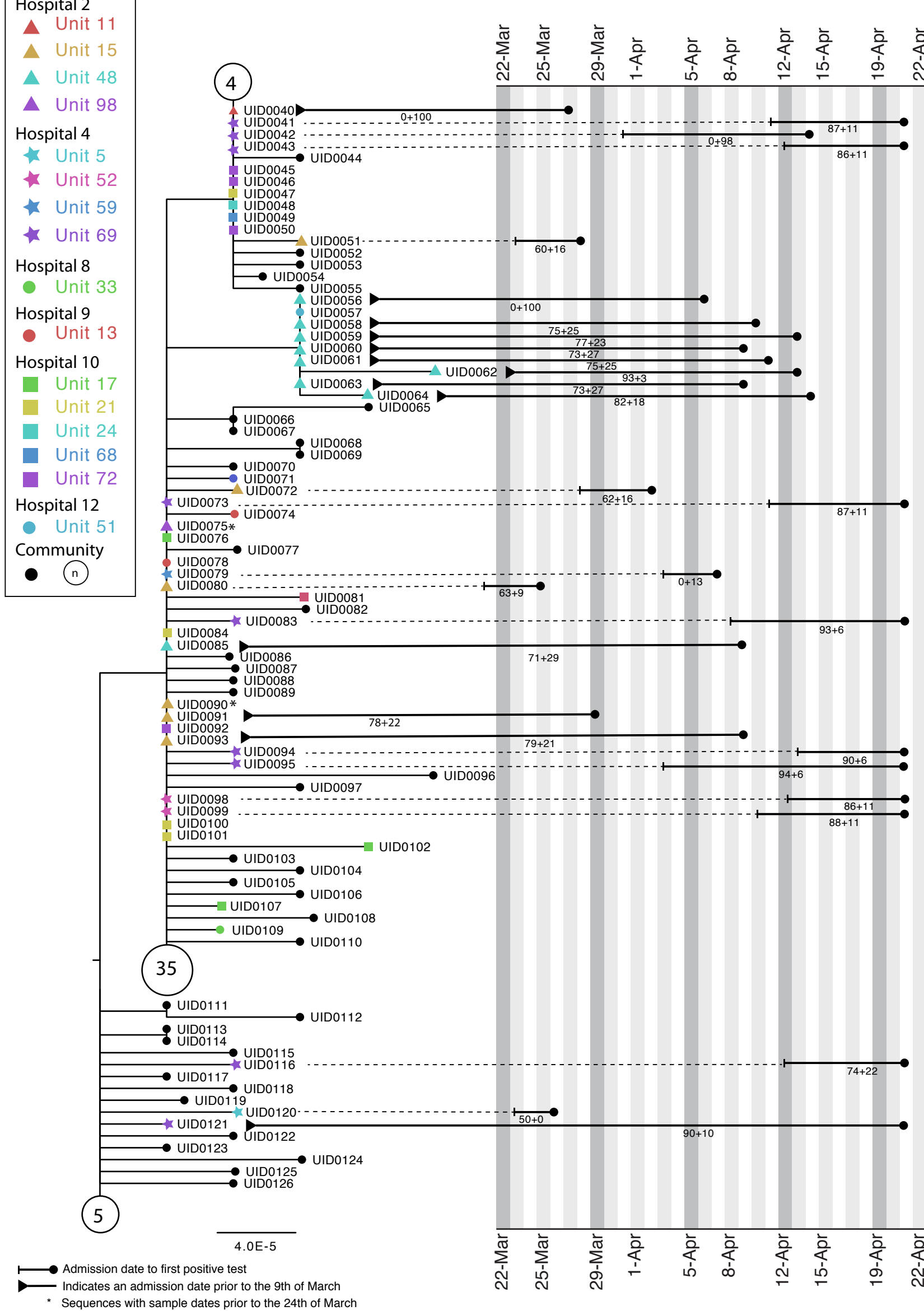


(b)



(a)**(b)**

- Hospital 1
 - Unit 89
- Hospital 2
 - Unit 11
 - Unit 15
 - Unit 48
 - Unit 98
- Hospital 4
 - Unit 5
 - Unit 52
 - Unit 59
 - Unit 69
- Hospital 8
 - Unit 33
- Hospital 9
 - Unit 13
- Hospital 10
 - Unit 17
 - Unit 21
 - Unit 24
 - Unit 68
 - Unit 72
- Hospital 12
 - Unit 51
- Community
 - (n)



Focus sample

UID0009

Report date	29-Oct-2020	Unit	Unit_93
Sample ID	-	Previous unit(s)	
Sample date	12-May-2020	Hospital	Hospital_5
COG-UK HOCI ID	-	Reporting hub	-
COG-UK ID	UID0009	Reported by	-
Admission date	21-Apr-2020	Symptomatic	Yes; onset date unknown

Report

Lineage: B.1.p73

Focus patient's sample sequence is closely matched to samples below, possibly linked by transmission.

⚠ Infection within unit is very highly probable* ⚠

Number	Sample ID	COG-UK ID	Other unit(s)	Sample date	Admission date	Type
1	-	UID0006	-	09-May-2020	30-Apr-2020	Patient
2	-	UID0018	-	09-May-2020	28-Apr-2020	Patient
3	-	UID0017	-	08-May-2020	01-May-2020	Patient
4	-	UID0022	-	12-May-2020	11-Apr-2020	Patient
5	-	UID0021	-	09-May-2020	01-May-2020	Patient
6	-	UID0032	-	05-May-2020	27-Apr-2020	Patient

Infection within hospital has low probability

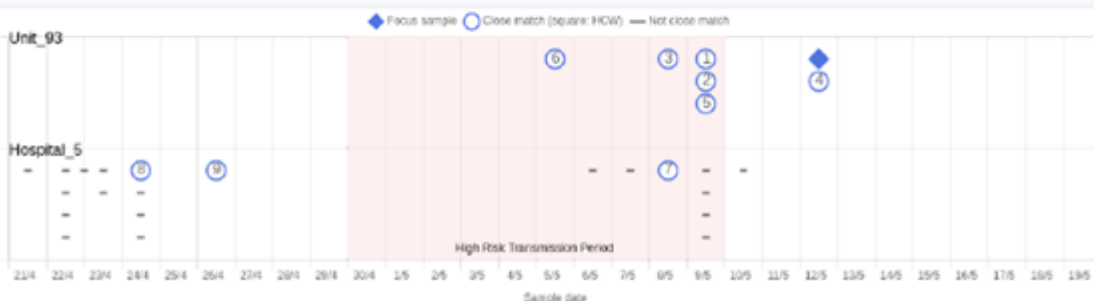
Number	Sample ID	COG-UK ID	Unit	Other unit(s)	Sample date	Admission date	Type
7	-	UID0025	Unit_92	-	08-May-2020	04-May-2020	Patient
8	-	UID0193	-	-	24-Apr-2020	-	Patient
9	-	UID0194	-	-	26-Apr-2020	-	Patient

Please check IPC data, and PATIENT and HCW movement, particularly for the 10-14 days preceding the date of the focus patient's sample.

- Infection from a visitor has low probability* (visitors not allowed on unit)
- Community-acquired infection has low probability*

* likelihood of transmission risk: 0-30% low; 30-50% moderately low; 50-70% probable; 70-85% high; 85-100% very high

Timeline



Focus sample

UID0025

Report date	29-Oct-2020	Unit	Unit_92
Sample ID	-	Previous unit(s)	
Sample date	08-May-2020	Hospital	Hospital_5
COG-UK HOCI ID	-	Reporting hub	-
COG-UK ID	UID0025	Reported by	-
Admission date	04-May-2020	Symptomatic	Yes; onset date unknown

Report

Lineage: B.1.p73

Focus patient's sample sequence is closely matched to samples below, possibly linked by transmission.

No matches from within unit

 Infection within hospital is probable 

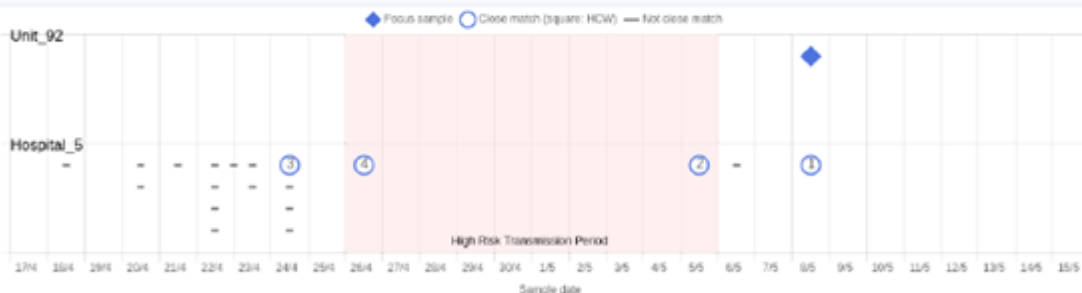
Number	Sample ID	COG-UK ID	Unit	Other unit(s)	Sample date	Admission date	Type
1	-	UID0017	Unit_93	-	08-May-2020	01-May-2020	Patient
2	-	UID0032	Unit_93	-	05-May-2020	27-Apr-2020	Patient
3	-	UID0193	-	-	24-Apr-2020	-	Patient
4	-	UID0194	-	-	26-Apr-2020	-	Patient

Please check IPC data, and PATIENT and HCW movement, particularly for the 10-14 days preceding the date of the focus patient's sample.

- Infection from a visitor has low probability* (visitors not allowed on unit)
- Community-acquired infection has moderately low probability*

* likelihood of transmission risk: 0-30% low; 30-50% moderately low; 50-70% probable; 70-85% high; 85-100% very high

Timeline



Generated on: 29-Oct-2020
GLUE version: 1.1.103

CoV-GLUE version: 0.1.13
COG-UK version: 0.1.6

HOCI version: 0.1.10
Author: Josh Singer <josh.singer@glasgow.ac.uk>