

1 **Genomic surveillance of *Clostridioides difficile* transmission and virulence in a healthcare setting**

2 Erin P. Newcomer<sup>1,2\*</sup>, Skye R. S. Fishbein<sup>1,3,\*</sup>, Kailun Zhang<sup>1,3</sup>, Tiffany Hink<sup>4</sup>, Kimberly A. Reske<sup>4</sup>,

3 Candice Cass<sup>4</sup>, Zainab H. Iqbal<sup>4</sup>, Emily L. Struttman<sup>4</sup>, Erik R. Dubberke<sup>4\*\*</sup>, Gautam Dantas<sup>1,2,3,5,6\*\*</sup>

4 <sup>1</sup>The Edison Family Center for Genome Sciences and Systems Biology, Washington University School of  
5 Medicine, St. Louis, Missouri, USA.

6 <sup>2</sup>Department of Biomedical Engineering, Washington University in St. Louis, St. Louis, Missouri, USA

7 <sup>3</sup>Department of Pathology and Immunology, Division of Laboratory and Genomic Medicine, Washington  
8 University School of Medicine, St. Louis, Missouri, USA.

9 <sup>4</sup>Division of Infectious Diseases; Washington University School of Medicine, St. Louis, Missouri, USA

10 <sup>5</sup>Department of Molecular Microbiology; Washington University School of Medicine, St. Louis, Missouri,  
11 USA

12 <sup>6</sup>Department of Pediatrics, Washington University School of Medicine, St. Louis, Missouri, USA

13 \*These authors contributed equally to this work

14 \*\*Corresponding authors: ERD: edubberk@wustl.edu; GD: dantas@wustl.edu

15 **Abstract**

16 *Clostridioides difficile* infection (CDI) is a major cause of healthcare-associated diarrhea,  
17 despite the widespread implementation of contact precautions for patients with CDI. Here, we  
18 investigate strain contamination in a hospital setting and genomic determinants of disease  
19 outcomes. Across two wards over six months, we selectively cultured *C. difficile* from patients  
20 (n=384) and their environments. Whole-genome sequencing (WGS) of 146 isolates revealed  
21 that most *C. difficile* isolates were from clade 1 (131/146, 89.7%), while only one isolate of the  
22 hypervirulent ST1 was recovered. Of culture-positive admissions, 17% of patients were  
23 diagnosed with CDI upon admission. We defined 29 strain networks at  $\leq 2$  core gene SNPs; 2 of  
24 these networks contain strains from different patients. Strain networks were temporally linked  
25 ( $p < 0.0001$ ). Across networks and over time, we found a minority of networks contained  
26 differences in phage populations. To understand genomic correlates of disease, we conducted  
27 WGS on an additional cohort of *C. difficile* (n=102 isolates) from the same hospital and  
28 confirmed that clade 1 isolates are responsible for most CDI cases. We found that while  
29 toxigenic *C. difficile* isolates are associated with the presence of *cdtR*, nontoxigenic isolates  
30 have an increased abundance of prophages. Our pangenomic analysis of clade 1 isolates  
31 suggests that while toxin genes (*tcdABER* and *cdtR*) were associated with CDI symptoms, they  
32 are dispensable for patient colonization. These data indicate toxigenic and nontoxigenic *C.*  
33 *difficile* contamination persists in a hospital setting and highlight further investigation into how  
34 accessory genomic repertoires contribute to *C. difficile* colonization and disease.

## 35 **Background**

36 *Clostridioides difficile* infection (CDI) is one of the most common healthcare-associated  
37 infections (HAIs) in the US and is the leading cause of healthcare-associated infectious  
38 diarrhea<sup>1,2</sup>. Since the early 2000s, *C. difficile* research has focused largely on hypervirulent  
39 strains, such as PCR ribotype 027<sup>1,3-6</sup>, which have been responsible for hospital-associated CDI  
40 epidemics. Strains of ribotype 027 were responsible for 51% and 84% of CDI cases in the US  
41 and Canada in 2005, respectively<sup>1,4,5</sup>. Since then, other circulating strains have emerged as the  
42 prevalent strains causative of CDI, such as 078 and 014/020<sup>7-9</sup>. One report indicated that the  
43 prevalence of PCR ribotype 027 decreased from 26.2% in 2012 to 16.9% in 2016<sup>9</sup>. As the  
44 landscape of *C. difficile* epidemiology continues to evolve, we must update our understanding of  
45 how various strains of this pathogen evolve, spread, and cause disease.

46 In addition to the changing prevalence of CDI-causing *C. difficile* strains, their  
47 transmission dynamics also appear to be evolving. In the late 1980s, it became clear that  
48 patients with active CDI shed spores onto their surroundings, leading to future CDI events in the  
49 healthcare setting<sup>1</sup>. Because of this, patients with active CDI are placed on contact precautions  
50 to prevent transmission to susceptible patients, which has been successful in reducing rates of  
51 CDI<sup>2,10</sup>. Nevertheless, while epidemiological estimates indicate that 20-42% of infections may be  
52 connected to a previous infection, multiple genomic studies fail to associate a CDI case to a  
53 previous case<sup>11-13</sup>. This suggests other potential sources of pathogen exposure in the hospital  
54 environment. While asymptomatic carriers of *C. difficile* have not been a significant focus of  
55 infection prevention efforts, studies have shown these carriers do shed viable, toxigenic *C.*  
56 *difficile* to their surroundings that could cause disease<sup>14</sup>. Several studies have shown evidence  
57 of a reduction in CDI cases if asymptomatic carriers are put on similar contact precautions to  
58 CDI patients<sup>15-17</sup>, but this has not been consistently found<sup>18</sup>. Correspondingly, it is critical to  
59 understand if *C. difficile* carriers are major contributors to new *C. difficile* acquisition or CDI  
60 manifestation in hospitalized patient populations.

61 *C. difficile* strains are categorized into five major clades and three additional cryptic  
62 clades. These clades encompass immense pangenomic diversity with many mobilizable  
63 chromosomal elements<sup>19,20</sup>, including numerous temperate phages that have potential  
64 influences over *C. difficile* toxin expression, sporulation, and metabolism<sup>21</sup>. Two major toxin loci,  
65 not required for viability, encode large multi-unit toxins that independently augment the virulence  
66 of *C. difficile*. Epithelial destruction and CDI have largely been attributed to the presence of  
67 pathogenicity locus (PaLoc) encoding toxins TcdA and TcdB. In addition, an accessory set of  
68 toxins (CdtA and CdtB) encoded at the binary toxin locus, may worsen disease symptoms<sup>22</sup>.

69 Yet, many nontoxigenic strains of *C. difficile* have been documented and are adept colonizers of  
70 the GI tract, even without the PaLoc<sup>23</sup>. As there has been continued debate about strain-specific  
71 virulence attributes<sup>24-26</sup>, it is important to investigate the extent of strain-level pangenomic  
72 diversity and consequences of such diversity on host disease<sup>27,28</sup>.

73 The purpose of this study was to evaluate the role of *C. difficile* strain diversity in  
74 colonization outcomes and hospital epidemiology. By sampling patients (n=384) and their  
75 environments for six months in two leukemia and hematopoietic stem cell (HCT) transplant  
76 wards at Barnes-Jewish Hospital in St. Louis, USA, we used isolate genomics to identify  
77 environmental contamination of both toxigenic (TCD) and nontoxigenic (NTCD) *C. difficile* by  
78 carriers and CDI patients, and corresponding transmission between both patient groups.  
79 Longitudinal strain tracking within these transmission networks revealed accessory gene flux of  
80 multi-drug resistance loci over the course of the study. Lastly, integration of isolate genomic  
81 data and CDI information from this prospective study with isolate genomic data from a  
82 complementary retrospective study of asymptomatic vs symptomatic *C. difficile* colonization in  
83 the same hospital<sup>29,30</sup> indicated that the clade 1 lineage, containing both toxigenic strains and  
84 nontoxigenic strains, dominates circulating populations of *C. difficile* in this hospital. Further, this  
85 lineage of *C. difficile* has significant variation in the PaLoc operon, and harbors other genetic  
86 factors that are associated with CDI symptoms in patients.

## 87 **Methods**

### 88 Study Design

89 This prospective observational study took place in the leukemia and hematopoietic stem cell  
90 transplant (HCT) wards at Barnes-Jewish Hospital (BJH) in St. Louis, Missouri, United States.  
91 Each ward consisted of two wings with 16 beds; on the acute leukemia ward we enrolled from  
92 both wings (32 beds) and on the HCT ward we enrolled on one wing (16 beds). The wards were  
93 sampled for 6 months from January 2019-July 2019 (acute leukemia) and 4 months from March  
94 2019-July 2019 (HCT). These units are located 2 floors apart in the same building.

### 95 Sample collection, selective culture, and isolate identification

96 Patients and their environments were sampled upon admission to a study ward and then weekly  
97 until discharge. Per hospital standards, bleach is used for daily and terminal discharge cleaning.  
98 From each patient, a stool specimen and/or rectal swab was collected as available. Remnant  
99 fecal samples from the BJH microbiology laboratory that were obtained during routine clinical  
100 care were also collected. Stool samples and rectal swabs collected on enrollment were  
101 refrigerated for up to 3 hours before processing. Specimens from all other timepoints were

102 stored in at -80°C in tryptic soy broth (TSB)/glycerol before processing. Environmental samples  
103 were collected from bedrails, keyboards, and sink surfaces using 3 E-swabs (Copan). If a  
104 surface was unable to be sampled, a swab was taken from the IV pump or nurse call button as  
105 an alternative. Swab eluate were stored at -80°C until processing.

106 Broth enrichment culture for *C. difficile* in Cycloserine Cefoxitin Mannitol Broth with  
107 Taurocholate and Lysozyme (CCMB-TAL) was performed on all admission specimens and  
108 checked for growth at 24h, 48h, and 7 days after inoculation. If that culture produced *C. difficile*,  
109 all other specimens collected from that patient and their surroundings were also cultured on  
110 Cycloserine-Cefoxitin Fructose Agar with Horse Blood and Taurocholate (CCFA-HT) agar.  
111 Colonies resembling *C. difficile* (large, spreading, grey, ground glass appearance) were picked  
112 by a trained microbiologist and sub-cultured onto a blood agar plate (BAP). Growth from the  
113 subculture plate was identified using Matrix-assisted laser desorption/ionization-time of flight  
114 mass spectrometry (MALDI-TOF MS). Upon identification, sweeps of *C. difficile* BAPs were  
115 collected in tryptic soy broth (TSB) and stored at -80C for sequencing. If both rectal swab  
116 sample and stool sample produced a *C. difficile* isolate, the stool isolate was preferentially used  
117 for analysis over the rectal swab isolate.

118 *C. difficile* toxin enzyme immunoassay (EIA) was conducted as part of routine clinical care  
119 based on clinical suspicion of CDI. To be diagnosed with *C. difficile* infection (CDI), a patient  
120 must have been EIA+ for *C. difficile* toxin (Alere TOX A/B II); those who weren't tested (due to  
121 no clinically significant diarrhea) or tested EIA- and were culture-positive for *C. difficile* were  
122 considered *C. difficile* carriers. Episodes of carriage or CDI are defined as the time from the first  
123 culture-positive specimen from a patient to the last culture-positive specimen during a given  
124 hospital admission.

125

#### 126 Short read sequencing and *de novo* genome assembly

127 Parameters used for computational tools are provided parenthetically. Total genomic DNA from  
128 *C. difficile* isolates was extracted from frozen plate scrapes using the QIAamp BiOstic  
129 Bacteremia DNA Kit (Qiagen) and quantified DNA with the PicoGreen dsDNA assay (Thermo  
130 Fisher Scientific). DNA from each isolate was diluted to a concentration of 0.5 ng/μL for library  
131 preparation using a modified Nextera kit (Illumina) protocol<sup>31</sup>. Sequencing libraries were pooled  
132 and sequenced on the NovaSeq 6000 platform (Illumina) to obtain 2 × 150 bp reads. Raw  
133 reads were demultiplexed by index pair and adapter sequencing trimmed and quality filtered  
134 using Trimmomatic (v0.38, SLIDINGWINDOW:4:20, LEADING:10, TRAILING:10,

135 MINLEN:60)<sup>32</sup>. Cleaned reads were assembled into draft genomes using Unicycler (v0.4.7)<sup>33</sup>.  
136 Draft genome quality was assessed using Quast<sup>34</sup>, BMAP<sup>35</sup>, and CheckM<sup>36</sup>, and genomes  
137 were accepted if they met the following quality standards: completeness greater than 90%,  
138 contamination less than 5%, N50 greater than 10,000 bp, and less than 500 contigs >1000bp.

#### 139 Isolate characterization and typing

140 A Mash Screen was used to identify likely related genomes from all NCBI reference genomes<sup>37</sup>.  
141 Average nucleotide identity(ANI) between the top three hits and the draft assembly was  
142 calculated using dnadiff<sup>38</sup>. Species were determined if an isolate had >75% alignment and  
143 >96% ANI<sup>39</sup> to a type strain, and were otherwise classified as genomospecies of the genus level  
144 taxonomy call.

145 In silico multilocus sequence typing (MLST) was determined for all *C. difficile* and  
146 genomospecies isolates using mlst<sup>40,41</sup>. Isolate contigs were annotated using Prokka<sup>42</sup> (v1.14.5,  
147 -mincontiglen 500, -force, -rnammer, -proteins GCF\_000210435.1\_ASM21043v1\_protein.faa<sup>43</sup>).  
148 *cdtAB* was determined to be a pseudogene if there were three hits to *cdtB*, indicating the  
149 damaged structure of the pseudogene<sup>44</sup>. *C. difficile* clade was determined using predefined  
150 clade-MLST relationships described in Knight, et al<sup>19</sup>.

#### 151 Phylogenetic analyses

152 The .gff files output by Prokka<sup>42</sup> were used as input for Panaroo (v1.2.10)<sup>45</sup> to construct a core  
153 genome alignment. The Panaroo alignment was used as input to construct a maximum-  
154 likelihood phylogenetic tree using Fasttree<sup>46</sup>. The output .newick file was visualized using the  
155 ggtree (v3.4.0)<sup>47</sup> package in R. Cryptic clade isolates were determined as such based on  
156 phylogenetic clustering with cryptic clade reference isolates.

#### 157 Core genome SNP analyses and network formation

158 We constructed a core gene alignment for each clade using Panaroo (v1.2.10) and calling  
159 MAFFT (v7.481). We then used Gubbins (v3.3.0) to identify recombination-filtered polymorphic  
160 sites, and constructed a recombination-free polymorphic site alignment using snp-sites (v2.4.0)  
161 25414349<sup>48</sup>. We finally extracted pairwise, recombination-filtered clade specific core-gene SNP  
162 distances using snp-dists (v0.8.2)(<https://github.com/tseemann/snp-dists>). Strain networks were  
163 determined by connecting isolates that were <=2 SNPs from one another.

#### 164 Phage identification and clustering

165 Isolate genomes were piped into Cenote-Taker 2<sup>49</sup> to identify contigs with end features as direct  
166 terminal repeats (DTRs) indicating circularity and inverted linear repeats (ITRs) or no features  
167 for linear sequences. Identified contigs were filtered by length and completeness to remove  
168 false positives. Length limits were 1,000 nucleotides (nt) for the detection of circularity, 4,000 nt  
169 for ITRs, and 5,000 nt for other linear sequences. The completeness was computed as a ratio  
170 between the length of our phage sequence and the length of matched reference genomes by  
171 CheckV<sup>50</sup> and the threshold was set to 10.0%. Phage contigs passing these two filters were  
172 then run through VIBRANT<sup>51</sup> with a “virome” flag to further remove obvious non-viral  
173 sequences<sup>51</sup>. Based on MIUViG recommended parameters<sup>52</sup>, phages were grouped into  
174 “populations” if they shared  $\geq 95\%$  nucleotide identity across  $\geq 85\%$  of the genome using  
175 BLASTN and a CheckV supporting code.

#### 176 Analysis of genotypic associations with disease severity

177 Two previously sequenced retrospective cohorts from the same hospital were included to  
178 increase statistical power<sup>29,53</sup>. In the analyses of toxigenic vs. nontoxigenic isolates from clade  
179 1, Pyseer<sup>54</sup> was run using a SNP distance matrix (using snp-dist as above), binary  
180 genotypes (presence or absence of *tcdB*), and Panaroo-derived gene presence/absence data.  
181 In the analysis of CDI suspicion, all isolates from clade 1 were used that represented one isolate  
182 per patient-episode. Isolates recovered from environmental surfaces were excluded. Using  
183 these assemblies, a core genome alignment was generated using Prokka<sup>42</sup> and Panaroo<sup>45</sup> as  
184 above. SNP distances were inferred from the core-gene alignment using snp-dists<sup>55</sup>. Binary  
185 phenotypes were coded for the variable CDI suspicion, whereby isolates associated with a  
186 clinically tested stool were associated with symptomatic colonization (TRUE). Isolates that were  
187 associated with a surveillance stool and had no clinical testing associated with that patient  
188 timepoint were coded as non-symptomatic colonization (FALSE). Gene candidates filtered  
189 based on ‘high-bse’, and were annotated HMMER on RefSeq databases and using a  
190 bacteriophage-specific tool VIBRANT<sup>51</sup>. Selected outputs were visualized in R using the beta  
191 coefficient as the x-axis and the  $-\log_{10}$ (likelihood ratio test p-value) as the y-axis.

#### 192 Reference assembly collection

193 We chose 23 reference assemblies from Knight, et al<sup>19</sup> for Figure 2c because of their MLST-  
194 clade associations (Supplementary Table 2). References span Clades 1-5 and cryptic clades C-  
195 1, C-2, and C-3, with one reference from each of the three most frequent MLSTs in each clade.  
196 Cryptic clade C-3 only had 2 reference assemblies available. References were annotated and  
197 included in phylogenetic tree construction as above.

198 All *Clostridioides difficile* genomes available on the National Institutes of Health (NIH) National  
199 Library of Medicine (NLM) were acquired for Figure 5c construction. References from NCBI  
200 (Supplementary Table 4) were included if they had less than 200 contigs. Assemblies that met  
201 these quality requirements were annotated and phylogenetically clustered as above.

## 202 **Results**

### 203 Surveillance of *C. difficile* reservoirs in hospital wards reveals patient colonization and 204 environmental contamination.

205 We prospectively collected patient and environmental samples to investigate genomic  
206 determinants of *C. difficile* carriage, transmission, and CDI (Figure 1). Across the study period,  
207 we enrolled 384 patients from 654 unique hospital admissions, and collected patient specimens  
208 upon admission and weekly thereafter (Supplementary Figure 1). We collected at least one  
209 specimen (clinical stool collected as part of routine care, study collected stool, or study collected  
210 rectal swab) from 364 admissions (94.8% of enrolled patients), for a total of 1244 patient  
211 specimens. We selectively cultured *C. difficile* from 43 rectal swabs and 108 stool samples, for a  
212 total of 151 culture-positive patient specimens. We also collected weekly swabs from the  
213 bedrails, sink surfaces, and in-room keyboards, for a total of 3045 swabs from each site. In total,  
214 22/398 (5.5%) of bedrail swabs cultured and 4/ 399 (1.0%) of keyboard swabs cultured were  
215 culture-positive for *C. difficile* (Figure 2a). *C. difficile* was never recovered from sink surfaces (all  
216 sinks on these units are hands-less activated) or other sampled sites. Collapsing multiple  
217 positive samples from the same patient admission results in 20 positive bedrails (20/79, 25.3%  
218 of all admissions with positive patient specimens) and 4 positive keyboards (4/79, 5.06% of all  
219 admissions with positive patient specimens) (Figure 2b).

### 220 *C. difficile* carriers outnumbered patients with CDI

221 Patients with CDI were identified through routine clinical care, with CDI defined as  
222 patients who had stool submitted for *C. difficile* testing, as ordered by the clinical team when  
223 suspicious for CDI, and who tested positive for *C. difficile* toxins by enzyme immunoassay  
224 (EIA+). Otherwise, if they were culture positive and EIA- or culture positive and not EIA tested,  
225 they were considered carriers. Results from selective culture indicated that 21.7% of unique  
226 admissions (79/364 admissions with available specimens) were culture-positive for *C. difficile* at  
227 some point during their admission (Figure 2b). Of culture-positive admissions, 17% (13/79) were  
228 EIA+ and diagnosed with CDI (13/364, 3.6% of all admissions with specimens available). The  
229 remaining 83% (66/79 admissions with specimens available) of culture-positive admissions



230 were termed carriers (Figure 2b). An additional nine admissions became EIA+ at some point  
231 during their stay for a total of 22 CDI cases, but seven did not have specimens available for  
232 culture and two were culture negative. The substantial detection of longitudinal patient *C.*  
233 *difficile* colonization prompted us to investigate the genomic correlates of *C. difficile*-associated  
234 disease and transmission in these two patient populations.

#### 235 Phylogenetic clustering reveals lack of hypervirulent strains, presence of cryptic clades

236 We conducted whole-genome sequencing to ascertain phylogenetic distances among isolates  
237 and to identify closely related strains of *C. difficile*. We identified 141 isolate genomes as *C.*  
238 *difficile* (using a 75% alignment and 96% average nucleotide identity [ANI] threshold). One  
239 isolate was identified as *Clostridium innocuum* and five isolates were classified as *C. difficile*  
240 genomospecies (92-93% ANI). To contextualize population structure, we applied a previously  
241 established MLST-derived clade definition to our isolate cohort<sup>19</sup>. The majority of *C. difficile*  
242 isolates were from Clade 1 (131/146, 89.7% of *C. difficile* and genomospecies, Figure 2c). Four  
243 patient-derived isolates were identified from clade 2, but only one was of the hypervirulent strain  
244 ST1 (PCR ribotype 027)<sup>6</sup>. We found that the distribution of STs associated with carriers was  
245 significantly different from that of STs associated with CDI patients ( $p < 0.001$ , Fisher's exact  
246 test) suggesting some strain-specificity to disease outcome.

247 Interestingly, the five genomospecies isolates clustered with other isolates belonging to  
248 a recently discovered *C. difficile* cryptic clade C-1 (Supplementary Figure 2). While cryptic  
249 clades are genomically divergent from *C. difficile*, these isolates can produce homologs to  
250 TcdA/B and cause CDI-like disease in humans<sup>19,56</sup>. In a clinical setting, they are frequently  
251 identified by MALDI-TOF MS as *C. difficile* and diagnosed as causative of CDI<sup>56</sup>. These data  
252 highlight the novel distribution of circulating *C. difficile* strains in the two study wards. While  
253 many patients with multiple isolates had homogeneous signatures of colonization (with closely  
254 related isolates), four patients (4/72, 6%) produced isolates from distinct ST types.

#### 255 Carriers and CDI patients contribute to transmission networks and environmental contamination

256 Given the predominance of Clade 1 isolates, we sought to identify clonal populations of *C.*  
257 *difficile* strains, indicative of direct *C. difficile* contamination (patient-environment) or  
258 transmission (patient-patient). We compared pairwise, recombination-filtered within-clade core  
259 gene single nucleotide polymorphism (SNP) distances to identify networks of transmission  
260 connecting isolates  $\leq 2$  SNPs apart (Supplementary Figure 4). We identified a total of 29 strain  
261 networks, 2 of which contain patient isolates from different patients (Figure 3a). These strain

262 networks were temporally linked, as there were significantly fewer days between same-network  
263 isolates than isolates from different networks ( $p < 2.2e-16$ , Wilcoxon, Figure 3b). We compared  
264 strain connections among a single patient's isolates from stool or rectal swab ('patient'), and  
265 between these isolates and environmental isolates from their immediate surroundings ('bedrail'  
266 or 'keyboard', Figure 3c). While the majority of bedrail isolates fell within the same network as  
267 patient isolates from that room (30 of 44 comparisons, 68%), 32% (14 of 44 comparisons) were  
268 genomically distinct, suggesting contamination from alternate sources. Keyboards were mostly  
269 colonized with distinct strains from the patient (22%, 2/9 comparisons were the same strain),  
270 indicating other routes of contamination ( $p < 0.05$ , Fisher's exact test, BH corrected. Figure 3c).  
271 Among the networks that contain multiple patients, we found no instances of potential  
272 transmission from the inhabitant of one room to the subsequent inhabitant. However, in both  
273 instances, each potential transmission was associated with a temporal overlap in patient stay in  
274 the same ward, providing epidemiological capacity for transmission ( $p < 0.05$ , Wilcoxon test).  
275 Importantly, we found no networks connecting patients with CDI to *C. difficile* carriers,  
276 suggesting successful containment through contact precaution protocols. These data highlight  
277 multiple sources of environmental contamination by *C. difficile* and prompted us to investigate  
278 the relationship between genetic factors and patient symptomology.

#### 279 Phage populations persist in circulating *C. difficile* networks

280 *C. difficile* isolates have an extensive pangenome, with genetic loci mobilized by  
281 conjugative elements and phages, and mobilizable elements playing a key role in *C. difficile*'s  
282 lifecycle<sup>57</sup>. Temperate phages, which can undergo lytic replication or insert into the host genome  
283 as a latent prophage, are the only phages that have been isolated for *C. difficile*<sup>58</sup>. To identify *C.*  
284 *difficile* prophage signatures and understand how dynamic they were in our strain networks, we  
285 analyzed our isolate genomes with Cenote-Taker 2 for putative phage contigs. After filtering for  
286 quality, we grouped contigs into phage populations (vOTUs) and quantified the alpha-diversity  
287 of phage populations in each isolate, and across MLST types (Figure 4a). ST42 and ST2, some  
288 of the most globally abundant ST types had the lowest diversity of phages in our cohort, though  
289 this negative correlation was not statistically significant across ST types (Figure 4b;  $R = -0.31$ ,  
290  $p = 0.12$ ). Our clonality-resolved strain networks allowed us to investigate phage flux over time.  
291 We found that the majority of networks (23/29) carried the same number of phages over time  
292 (Figure 4c), suggesting persistent roles in *C. difficile* biology. Interestingly, we found that  
293 nontoxigenic isolates had a higher diversity of phage populations relative to toxigenic isolates

294 (Figure 4d). These data suggest distinct selective pressures on temperate phages in *C. difficile*  
295 related to toxin gene presence.

#### 296 Accessory genomic elements are associated with host CDI symptoms

297 Despite evidence of transmission in this prospective study, a minority of patients were  
298 diagnosed with CDI relative to those asymptomatically colonized with *C. difficile* in part due to  
299 the presence of nontoxigenic *C. difficile* isolates (Figure 2b). To power our investigation of  
300 virulence determinants across patient-colonizing *C. difficile* strains, we performed whole  
301 genome sequencing on 102 additional patient-derived *C. difficile* isolates from a previously  
302 described *C. difficile*-colonized/CDI cohort from the same hospital<sup>29</sup>, where all patients had  
303 clinical suspicion of CDI (CDI suspicion), defined by a clinician ordering an EIA test during  
304 patient admission. Using an MLST-based clade definition as above, we identified that most CDI  
305 cases result from isolates within clade 1, though clade 2 isolates were more likely to be  
306 associated with CDI status (Figure 5a). The latter finding supports previous data indicating that  
307 clade 2 isolates are hypervirulent, often attributed to the presence of the binary toxin operon or  
308 increased expression from the PaLoc<sup>22,59,60</sup>. Meanwhile, some clade 1 isolates contain no  
309 toxins, indicating a diversity of colonization strategies in this lineage. Pangenomic comparison of  
310 nontoxigenic versus toxigenic isolates revealed that in addition to the PaLoc, the majority of our  
311 toxigenic isolates from clade 1 (95/131 of our cohort) possess remnants of the binary toxin  
312 operon (Figure 5b, *cdtR* and *cdtA/B* pseudogenes). Given the previous report that full-length  
313 *cdtAB* was identified only within Clades 2, 3, and 5<sup>19</sup>, we investigated the conservation of *cdtR*  
314 (the transcriptional regulator of the binary toxin locus) across *C. difficile* strains (containing 5  
315 lineages). We additionally examined >1400 *C. difficile* genome assemblies from NCBI  
316 (Supplementary Table 4, Figure 5c). *cdtR* (unlike *cdtAB*) was dispersed across clade 1 and  
317 significantly associated with *tcdB* (Figure 5d, Fisher's exact test, BH corrected), suggesting a  
318 selective pressure to maintain some element of both toxin loci in these isolates. Notably, these  
319 operons are not syntenic, further underlining the significance of the association. From this  
320 association, we sought to further understand why some toxigenic clade 1 isolates cause CDI  
321 and some colonize without symptoms. Using 148 toxigenic clade 1 isolates collected from this  
322 study and two previous studies from the same hospital<sup>29,53</sup>, we utilized a bacterial GWAS  
323 approach, *pyseer*<sup>54</sup>, that identifies genetic traits associated with strains corresponding to  
324 patients with CDI symptoms. Using CDI suspicion (see Methods) as an outcome variable, we  
325 found that, multiple amidases (including *cwiD*), putative transcriptional regulators, and many  
326 genes of unknown function were enriched in isolates associated with CDI symptoms (Figure

327 5e). These data indicate that the most prevalent, circulating *Cd* strains that cause CDI are not  
328 the hypervirulent clade 2 strains, but highlight the possibility that remnant genomic features from  
329 epidemic strains and other features may contribute to virulence in this hospital clade of *C.*  
330 *difficile*.

### 331 **Discussion**

332 Through our prospective genomics study of two hospital wards, we were able to identify  
333 connections between the contamination of different surfaces and the strains carried by  
334 hospitalized patients and quantify some spread between carriers. Our estimates of the  
335 prevalence of patients with CDI (3.8%) agree with other estimates of 2-4% CDI in patients with  
336 cancer<sup>61-63</sup>. While many studies have quantified surface contamination, few have had the  
337 genomic resolution to identify clonality between isolates indicating transmission or patient  
338 shedding<sup>64-66</sup>. We observed distinct patterns of contamination between a patient's bedrail and  
339 the corresponding room keyboard, supporting the notion that the bedrail could be one of  
340 multiple critical points of transmission in a hospital setting. Further, we did not identify any  
341 instances of CDI that could be genomically linked to an earlier CDI case or *C. difficile* carrier.  
342 Despite the small sample size, these data support the continued use of contact precautions for  
343 CDI patients<sup>18</sup>.

344 Our data suggests the need to continually update our understanding of CDI-causing *C.*  
345 *difficile* strains beyond previous epidemic strains to clarify mechanisms of how the most  
346 prevalent strains relate to transmission and disease. Across 146 patient specimens, we only  
347 identified one incidence of the epidemic ST1 strain. This ribotype caused one case of CDI within  
348 our cohort, corroborating the decline in this epidemic lineage<sup>67</sup>. Because the overall burden of  
349 Clade 1 isolates was so high, we hypothesize that its ability to colonize without causing CDI  
350 could allow for a substantial expansion of transmission networks (especially for the case of  
351 nontoxigenic strains). While Clade 1 isolates associated with CDI symptoms are expectedly  
352 toxigenic (containing the toxin genes in the PaLoc), we also found an enrichment in two different  
353 amidase genes, that could either contribute to differences in germination rate or possess  
354 endolysin function<sup>68 69</sup>. How the function of such a gene contributes to an increase in  
355 symptomology remains to be understood. Further, we confirmed a genetic relationship between  
356 *cdtR* and *tcdB* across *C. difficile* lineages that indicates some evolutionary pressure for  
357 maintaining the regulatory gene of the less prevalent toxin operon (*cdtR*). This phylogenomic  
358 analysis supports recent functional data from clade 2 isolates, where the presence of full-length  
359 *cdtR* increases the expression of *tcdB* and disease severity in an animal model of CDI<sup>57</sup>. While

360 this was previously suggested *in vitro*, it is unclear how generalizable this relationship is across  
361 lineages<sup>59</sup>. In fact, we predict that clade 1 isolates containing only *cdtR* and the PaLoc may  
362 produce more toxin *in vivo*. Future studies are warranted to investigate the role of both classes  
363 of genes implicated in this phenotype.

364 Our study contextualizes the need for investigating *C. difficile* evolution within patients  
365 over time, especially concerning functional mobile units such as temperate phages. We  
366 examined phage populations in our isolates as they are a relevant mobile unit of the *C. difficile*  
367 pangenome and their stability over time has not been systematically investigated. While we find  
368 that the majority of *C. difficile* strains maintain their diversity of phage populations over time, we  
369 acknowledge that hospital admission is a prescribed period of time and we may be  
370 underestimating the amount that phage diversity changes in isolates over longer periods of time  
371 *in vivo*. Our quantification of increased phage diversity in nontoxicogenic isolates suggests phage  
372 niche specialization based on the presence of the PaLoc. It is noteworthy that early  
373 characterization of the PaLoc operon indicated that it was integrated into the *C. difficile*  
374 chromosome by an ancient prophage<sup>70-72</sup>. Future work is required to understand how persistent  
375 phages function during *C. difficile* growth and pathogenesis<sup>58</sup>.

376 Our study has a number of important limitations. As this study focused on *C. difficile*  
377 colonization, disease, and transmission in two wards in the same hospital, studies with  
378 increased sample size or meta-analysis studies are necessary to understand generalizable  
379 epidemiological measurements of *C. difficile*-patient dynamics<sup>73</sup>. Additionally, our study protocol  
380 allowed for culturing all environmental/patient specimens from a carrier or patient with CDI.  
381 Thus, it is possible that our estimate of carriage in this study population is an overestimate.  
382 Finally, we note the evidence for multi-strain colonization within a single patient (Patient 2330).  
383 Given our approach of only culturing and sequencing single isolates per patient timepoint, future  
384 studies are needed to investigate the extent of within-patient *C. difficile* strain diversity by  
385 interrogating additional cultured isolates per samples<sup>74</sup> or via metagenomic methods.

386 Despite these limitations, this work allows us to understand an updated genomic picture  
387 of circulating *C. difficile* in hospital-associated patients: how strains spread, their evolution, and  
388 their virulence potential in this study population. Indeed, though much human and animal  
389 research has focused on epidemic strains that are two decades old, we and others have  
390 identified more disease and colonization from distinct lineages of *C. difficile*, namely clade 1  
391 lineages. Moreover, within this lineage we found a mosaic representation of the PaLoc that  
392 highlighted the possibility of different mechanisms of colonization and virulence by this

393 population of *C. difficile*. Future studies utilizing other human cohorts or animal models are  
394 warranted to investigate disease and pathogenicity caused by Clade 1 *C. difficile* strains.

395

#### 396 **List of Abbreviations**

397 BAP: blood agar plate

398 CCFA-HT: Cycloserine-Cefoxitin Fructose Agar with Horse Blood and Taurocholate

399 CCMB-TAL: Cycloserine Cefoxitin Mannitol Broth with Taurocholate and Lysozyme

400 CDI: *Clostridioides difficile* infection

401 EIA: enzyme immunoassay

402 HAI: healthcare-associated infection

403 HGT: horizontal gene transfer

404 MALDI-TOF MS: Matrix-assisted laser desorption/ionization-time of flight mass spectrometry

405 NTCD: non-toxigenic *C. difficile*

406 PaLoc: pathogenicity locus

407 TCD: toxigenic *C. difficile*

408 TSB: tryptic soy broth

409

#### 410 **Declarations**

##### 411 **Ethics approval and consent to participate**

412 The study protocol was approved by the Washington University Human Research Protection  
413 Office (IRB #201810103). All participants provided written informed consent.

##### 414 **Consent for publication**

415 Not applicable.

##### 416 **Availability of data and materials**

417 The datasets generated and analyzed during the current study are available in NCBI GenBank  
418 under BioProject accession no. PRJNA980715.

##### 419 **Competing interests**

420 The authors declare that they have no competing interests.

##### 421 **Funding**

422 This work was supported in part by an award to ERD and GD through the Foundation for  
423 Barnes-Jewish Hospital and Institute of Clinical and Translational Sciences. This publication  
424 was supported by the NIH/National Center for Advancing Translational Sciences (NCATS),  
425 grant UL1 TR002345 (PI: B. Evanoff). This work was also supported by funding through the  
426 CDC BAA #200-2018-02926 under PI Erik Dubberke. SRSF is supported by the National  
427 Institute of Child Health and Human Development (NICHD: <https://www.nicdhd.nih.gov>) of the  
428 NIH under award number T32 HD004010 (PI: P. Tarr). The conclusions from this study  
429 represent those of the authors and do not represent positions of the funding agencies.

#### 430 **Authors' contributions**

431 SRSF, KAR, ERD, and GD participated in idea formulation and funding for this project. TH,  
432 KAR, CC, ZHI, ELS, and ERD conducted participant enrollment, sample collection, and  
433 microbiological isolation. EPN, SRSF, KZ, and GD conducted all sequencing analysis and figure  
434 generation. EPN and SRSF completed the writing of the manuscript. All authors read and  
435 approved the final manuscript.

#### 436 **Acknowledgements**

437 The authors are grateful for members of the Dantas lab for their helpful feedback on the data  
438 analysis and preparation of the manuscript. The authors would also like to thank the Edison  
439 Family Center for Genome Sciences and Systems Biology staff, Eric Martin, Brian Koebbe,  
440 MariaLynn Crosby, and Jessica Hoisington-López for their expertise and support in  
441 sequencing/data analysis.

## 442 References

- 443 1 Czepiel, J. *et al.* Clostridium difficile infection: review. *Eur J Clin Microbiol Infect Dis* **38**,  
444 1211-1221, doi:10.1007/s10096-019-03539-6 (2019).
- 445 2 McDonald, L. C. *et al.* Clinical Practice Guidelines for Clostridium difficile Infection in  
446 Adults and Children: 2017 Update by the Infectious Diseases Society of America (IDSA)  
447 and Society for Healthcare Epidemiology of America (SHEA). *Clin Infect Dis* **66**, e1-e48,  
448 doi:10.1093/cid/cix1085 (2018).
- 449 3 Clements, A. C., Magalhaes, R. J., Tatem, A. J., Paterson, D. L. & Riley, T. V. Clostridium  
450 difficile PCR ribotype 027: assessing the risks of further worldwide spread. *Lancet Infect*  
451 *Dis* **10**, 395-404, doi:10.1016/S1473-3099(10)70080-3 (2010).
- 452 4 McDonald, L. C. *et al.* An epidemic, toxin gene-variant strain of Clostridium difficile. *N*  
453 *Engl J Med* **353**, 2433-2441, doi:10.1056/NEJMoa051590 (2005).
- 454 5 Loo, V. G. *et al.* A predominantly clonal multi-institutional outbreak of Clostridium  
455 difficile-associated diarrhea with high morbidity and mortality. *N Engl J Med* **353**, 2442-  
456 2449, doi:10.1056/NEJMoa051639 (2005).
- 457 6 Fatima, R. & Aziz, M. The Hypervirulent Strain of Clostridium Difficile: NAP1/B1/027 - A  
458 Brief Overview. *Cureus* **11**, e3977, doi:10.7759/cureus.3977 (2019).
- 459 7 Goorhuis, A. *et al.* Emergence of Clostridium difficile infection due to a new  
460 hypervirulent strain, polymerase chain reaction ribotype 078. *Clin Infect Dis* **47**, 1162-  
461 1170, doi:10.1086/592257 (2008).
- 462 8 Bauer, M. P. *et al.* Clostridium difficile infection in Europe: a hospital-based survey.  
463 *Lancet* **377**, 63-73, doi:10.1016/S0140-6736(10)61266-4 (2011).
- 464 9 Giancola, S. E., Williams, R. J., 2nd & Gentry, C. A. Prevalence of the Clostridium difficile  
465 BI/NAP1/027 strain across the United States Veterans Health Administration. *Clin*  
466 *Microbiol Infect* **24**, 877-881, doi:10.1016/j.cmi.2017.11.011 (2018).
- 467 10 Balsells, E. *et al.* Infection prevention and control of Clostridium difficile: a global review  
468 of guidelines, strategies, and recommendations. *J Glob Health* **6**, 020410,  
469 doi:10.7189/jogh.06.020410 (2016).
- 470 11 Kong, L. Y. *et al.* Clostridium difficile: Investigating Transmission Patterns Between  
471 Infected and Colonized Patients Using Whole Genome Sequencing. *Clin Infect Dis* **68**,  
472 204-209, doi:10.1093/cid/ciy457 (2019).
- 473 12 Svenungsson, B. *et al.* Epidemiology and molecular characterization of Clostridium  
474 difficile strains from patients with diarrhea: low disease incidence and evidence of  
475 limited cross-infection in a Swedish teaching hospital. *J Clin Microbiol* **41**, 4031-4037,  
476 doi:10.1128/JCM.41.9.4031-4037.2003 (2003).
- 477 13 Durham, D. P., Olsen, M. A., Dubberke, E. R., Galvani, A. P. & Townsend, J. P. Quantifying  
478 Transmission of Clostridium difficile within and outside Healthcare Settings. *Emerg*  
479 *Infect Dis* **22**, 608-616, doi:10.3201/eid2204.150455 (2016).
- 480 14 Warren, B. G. *et al.* The Impact of Infection Versus Colonization on Clostridioides difficile  
481 Environmental Contamination in Hospitalized Patients With Diarrhea. *Open Forum Infect*  
482 *Dis* **9**, ofac069, doi:10.1093/ofid/ofac069 (2022).



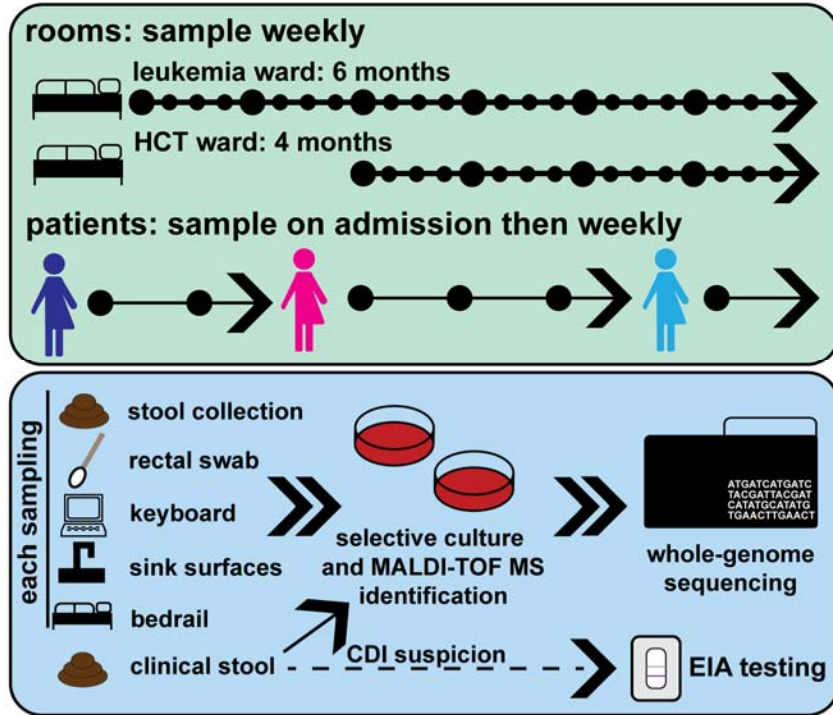
- 483 15 Longtin, Y. *et al.* Effect of Detecting and Isolating *Clostridium difficile* Carriers at Hospital  
484 Admission on the Incidence of *C. difficile* Infections: A Quasi-Experimental Controlled  
485 Study. *JAMA Intern Med* **176**, 796-804, doi:10.1001/jamainternmed.2016.0177 (2016).
- 486 16 Xiao, Y. *et al.* Impact of Isolating *Clostridium difficile* Carriers on the Burden of Isolation  
487 Precautions: A Time Series Analysis. *Clin Infect Dis* **66**, 1377-1382,  
488 doi:10.1093/cid/cix1024 (2018).
- 489 17 Grigoras, C. A., Zervou, F. N., Zacharioudakis, I. M., Siettos, C. I. & Mylonakis, E. Isolation  
490 of *C. difficile* Carriers Alone and as Part of a Bundle Approach for the Prevention of  
491 *Clostridium difficile* Infection (CDI): A Mathematical Model Based on Clinical Study Data.  
492 *PLoS One* **11**, e0156577, doi:10.1371/journal.pone.0156577 (2016).
- 493 18 Morgan, D. J. *et al.* The Impact of Universal Glove and Gown Use on *Clostridioides*  
494 *Difficile* Acquisition: A Cluster-Randomized Trial. *Clin Infect Dis* **76**, e1202-e1207,  
495 doi:10.1093/cid/ciac519 (2023).
- 496 19 Knight, D. R. *et al.* Major genetic discontinuity and novel toxigenic species in  
497 *Clostridioides difficile* taxonomy. *Elife* **10**, doi:10.7554/eLife.64325 (2021).
- 498 20 Mullany, P., Allan, E. & Roberts, A. P. Mobile genetic elements in *Clostridium difficile*  
499 and their role in genome function. *Res Microbiol* **166**, 361-367,  
500 doi:10.1016/j.resmic.2014.12.005 (2015).
- 501 21 Fortier, L. C. Bacteriophages Contribute to Shaping *Clostridioides (Clostridium) difficile*  
502 Species. *Front Microbiol* **9**, 2033, doi:10.3389/fmicb.2018.02033 (2018).
- 503 22 Gerding, D. N., Johnson, S., Rupnik, M. & Aktories, K. *Clostridium difficile* binary toxin  
504 CDT: mechanism, epidemiology, and potential clinical importance. *Gut Microbes* **5**, 15-  
505 27, doi:10.4161/gmic.26854 (2014).
- 506 23 Gerding, D. N. *et al.* Administration of spores of nontoxigenic *Clostridium difficile* strain  
507 M3 for prevention of recurrent *C. difficile* infection: a randomized clinical trial. *JAMA*  
508 **313**, 1719-1727, doi:10.1001/jama.2015.3725 (2015).
- 509 24 Carlson, P. E., Jr. *et al.* The relationship between phenotype, ribotype, and clinical  
510 disease in human *Clostridium difficile* isolates. *Anaerobe* **24**, 109-116,  
511 doi:10.1016/j.anaerobe.2013.04.003 (2013).
- 512 25 Walk, S. T. *et al.* *Clostridium difficile* ribotype does not predict severe infection. *Clin*  
513 *Infect Dis* **55**, 1661-1668, doi:10.1093/cid/cis786 (2012).
- 514 26 Aitken, S. L. *et al.* In the Endemic Setting, *Clostridium difficile* Ribotype 027 Is Virulent  
515 But Not Hypervirulent. *Infect Control Hosp Epidemiol* **36**, 1318-1323,  
516 doi:10.1017/ice.2015.187 (2015).
- 517 27 Pettit, L. J. *et al.* Functional genomics reveals that *Clostridium difficile* Spo0A  
518 coordinates sporulation, virulence and metabolism. *BMC Genomics* **15**, 160,  
519 doi:10.1186/1471-2164-15-160 (2014).
- 520 28 Awad, M. M., Johanesen, P. A., Carter, G. P., Rose, E. & Lyras, D. *Clostridium difficile*  
521 virulence factors: Insights into an anaerobic spore-forming pathogen. *Gut Microbes* **5**,  
522 579-593, doi:10.4161/19490976.2014.969632 (2014).
- 523 29 Fishbein, S. R. *et al.* Multi-omics investigation of *Clostridioides difficile*-colonized  
524 patients reveals pathogen and commensal correlates of *C. difficile* pathogenesis. *Elife*  
525 **11**, doi:10.7554/eLife.72801 (2022).

- 526 30 Dubberke, E. R. *et al.* Clostridium difficile colonization among patients with clinically  
527 significant diarrhea and no identifiable cause of diarrhea. *Infect Control Hosp Epidemiol*  
528 **39**, 1330-1333, doi:10.1017/ice.2018.225 (2018).
- 529 31 Baym, M. *et al.* Inexpensive multiplexed library preparation for megabase-sized  
530 genomes. *PLoS One* **10**, e0128036, doi:10.1371/journal.pone.0128036 (2015).
- 531 32 Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina  
532 sequence data. *Bioinformatics* **30**, 2114-2120, doi:10.1093/bioinformatics/btu170  
533 (2014).
- 534 33 Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial genome  
535 assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595,  
536 doi:10.1371/journal.pcbi.1005595 (2017).
- 537 34 Gurevich, A., Saveliev, V., Vyahhi, N. & Tesler, G. QUAST: quality assessment tool for  
538 genome assemblies. *Bioinformatics* **29**, 1072-1075, doi:10.1093/bioinformatics/btt086  
539 (2013).
- 540 35 Bushnell, B. *BBMap: A Fast, Accurate, Splice-Aware Aligner*,  
541 <<https://www.osti.gov/servlets/purl/1241166>> (2014).
- 542 36 Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:  
543 assessing the quality of microbial genomes recovered from isolates, single cells, and  
544 metagenomes. *Genome Res* **25**, 1043-1055, doi:10.1101/gr.186072.114 (2015).
- 545 37 Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using  
546 MinHash. *Genome Biol* **17**, 132, doi:10.1186/s13059-016-0997-x (2016).
- 547 38 Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol*  
548 **5**, R12, doi:10.1186/gb-2004-5-2-r12 (2004).
- 549 39 Richter, M. & Rossello-Mora, R. Shifting the genomic gold standard for the prokaryotic  
550 species definition. *Proc Natl Acad Sci U S A* **106**, 19126-19131,  
551 doi:10.1073/pnas.0906412106 (2009).
- 552 40 Seemann, T. *mlst*, <<https://github.com/tseemann/mlst>> (  
553 41 Jolley, K. A. & Maiden, M. C. BIGSdb: Scalable analysis of bacterial genome variation at  
554 the population level. *BMC Bioinformatics* **11**, 595, doi:10.1186/1471-2105-11-595  
555 (2010).
- 556 42 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068-  
557 2069, doi:10.1093/bioinformatics/btu153 (2014).
- 558 43 He, M. *et al.* Evolutionary dynamics of Clostridium difficile over short and long time  
559 scales. *Proc Natl Acad Sci U S A* **107**, 7527-7532, doi:10.1073/pnas.0914322107 (2010).
- 560 44 Carter, G. P. *et al.* Binary toxin production in Clostridium difficile is regulated by CdtR, a  
561 LytTR family response regulator. *J Bacteriol* **189**, 7290-7301, doi:10.1128/JB.00731-07  
562 (2007).
- 563 45 Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo  
564 pipeline. *Genome Biol* **21**, 180, doi:10.1186/s13059-020-02090-4 (2020).
- 565 46 Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree: computing large minimum evolution  
566 trees with profiles instead of a distance matrix. *Mol Biol Evol* **26**, 1641-1650,  
567 doi:10.1093/molbev/msp077 (2009).
- 568 47 Yu, G. Using ggtree to Visualize Data on Tree-Like Structures. *Curr Protoc Bioinformatics*  
569 **69**, e96, doi:10.1002/cpbi.96 (2020).

- 570 48 Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA  
571 alignments. *Microb Genom* **2**, e000056, doi:10.1099/mgen.0.000056 (2016).
- 572 49 Tisza, M. J., Belford, A. K., Dominguez-Huerta, G., Bolduc, B. & Buck, C. B. Cenote-Taker 2  
573 democratizes virus discovery and sequence annotation. *Virus Evol* **7**, veaa100,  
574 doi:10.1093/ve/veaa100 (2021).
- 575 50 Nayfach, S. *et al.* CheckV assesses the quality and completeness of metagenome-  
576 assembled viral genomes. *Nat Biotechnol* **39**, 578-585, doi:10.1038/s41587-020-00774-7  
577 (2021).
- 578 51 Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and  
579 curation of microbial viruses, and evaluation of viral community function from genomic  
580 sequences. *Microbiome* **8**, 90, doi:10.1186/s40168-020-00867-0 (2020).
- 581 52 Roux, S. *et al.* Minimum Information about an Uncultivated Virus Genome (MIUViG). *Nat*  
582 *Biotechnol* **37**, 29-37, doi:10.1038/nbt.4306 (2019).
- 583 53 Fishbein, S. R. S. *et al.* Randomized Controlled Trial of Oral Vancomycin Treatment in  
584 Clostridioides difficile-Colonized Patients. *mSphere* **6**, doi:10.1128/mSphere.00936-20  
585 (2021).
- 586 54 Lees, J. A., Galardini, M., Bentley, S. D., Weiser, J. N. & Corander, J. pyseer: a  
587 comprehensive tool for microbial pangenome-wide association studies. *Bioinformatics*  
588 **34**, 4310-4312, doi:10.1093/bioinformatics/bty539 (2018).
- 589 55 Seemann, T. *snippy: fast bacterial variant calling from NGS reads*,  
590 <<https://github.com/tseemann/snippy>> (2015).
- 591 56 Williamson, C. H. D. *et al.* Identification of novel, cryptic Clostridioides species isolates  
592 from environmental samples collected from diverse geographical locations. *Microb*  
593 *Genom* **8**, doi:10.1099/mgen.0.000742 (2022).
- 594 57 Dong, Q. *et al.* Virulence and genomic diversity among clinical isolates of ST1  
595 (BI/NAP1/027) Clostridioides difficile. *Cell Rep* **42**, 112861,  
596 doi:10.1016/j.celrep.2023.112861 (2023).
- 597 58 Heuler, J., Fortier, L. C. & Sun, X. Clostridioides difficile phage biology and application.  
598 *FEMS Microbiol Rev* **45**, doi:10.1093/femsre/fuab012 (2021).
- 599 59 Lyon, S. A., Hutton, M. L., Rood, J. I., Cheung, J. K. & Lyras, D. CdtR Regulates TcdA and  
600 TcdB Production in Clostridium difficile. *PLoS Pathog* **12**, e1005758,  
601 doi:10.1371/journal.ppat.1005758 (2016).
- 602 60 Dong, Q. *et al.* Virulence and genomic diversity among clinical isolates of ST1  
603 (BI/NAP1/027) Clostridioides difficile. *bioRxiv*, doi:10.1101/2023.01.12.523823 (2023).
- 604 61 Zheng, Y. *et al.* Clostridium difficile colonization in preoperative colorectal cancer  
605 patients. *Oncotarget* **8**, 11877-11886, doi:10.18632/oncotarget.14424 (2017).
- 606 62 Jain, T. *et al.* Clostridium Difficile Colonization in Hematopoietic Stem Cell Transplant  
607 Recipients: A Prospective Study of the Epidemiology and Outcomes Involving Toxigenic  
608 and Nontoxigenic Strains. *Biol Blood Marrow Transplant* **22**, 157-163,  
609 doi:10.1016/j.bbmt.2015.07.020 (2016).
- 610 63 Kamboj, M., Gennarelli, R. L., Brite, J., Sepkowitz, K. & Lipitz-Snyderman, A. Risk for  
611 Clostridioides difficile Infection among Older Adults with Cancer. *Emerg Infect Dis* **25**,  
612 1683-1689, doi:10.3201/eid2509.181142 (2019).

- 613 64 Claro, T., Daniels, S. & Humphreys, H. Detecting *Clostridium difficile* spores from  
614 inanimate surfaces of the hospital environment: which method is best? *J Clin Microbiol*  
615 **52**, 3426-3428, doi:10.1128/JCM.01011-14 (2014).
- 616 65 Kumar, N. *et al.* Genome-Based Infection Tracking Reveals Dynamics of *Clostridium*  
617 *difficile* Transmission and Disease Recurrence. *Clin Infect Dis* **62**, 746-752,  
618 doi:10.1093/cid/civ1031 (2016).
- 619 66 Kiersnowska, Z. M., Lemiech-Mirowska, E., Michalkiewicz, M., Sierocka, A. & Marczak,  
620 M. Detection and Analysis of *Clostridioides difficile* Spores in a Hospital Environment. *Int*  
621 *J Environ Res Public Health* **19**, doi:10.3390/ijerph192315670 (2022).
- 622 67 Snyderman, D. R. *et al.* Epidemiologic trends in *Clostridioides difficile* isolate ribotypes in  
623 United States from 2011 to 2016. *Anaerobe* **63**, 102185,  
624 doi:10.1016/j.anaerobe.2020.102185 (2020).
- 625 68 Diaz, O. R., Sayer, C. V., Popham, D. L. & Shen, A. *Clostridium difficile* Lipoprotein GerS Is  
626 Required for Cortex Modification and Thus Spore Germination. *mSphere* **3**,  
627 doi:10.1128/mSphere.00205-18 (2018).
- 628 69 Wydau-Dematteis, S. *et al.* Cwp19 Is a Novel Lytic Transglycosylase Involved in  
629 Stationary-Phase Autolysis Resulting in Toxin Release in *Clostridium difficile*. *mBio* **9**,  
630 doi:10.1128/mBio.00648-18 (2018).
- 631 70 Canchaya, C., Proux, C., Fournous, G., Bruttin, A. & Brussow, H. Prophage genomics.  
632 *Microbiol Mol Biol Rev* **67**, 238-276, table of contents, doi:10.1128/MMBR.67.2.238-  
633 276.2003 (2003).
- 634 71 Proux, C. *et al.* The dilemma of phage taxonomy illustrated by comparative genomics of  
635 Sfi21-like Siphoviridae in lactic acid bacteria. *J Bacteriol* **184**, 6026-6036,  
636 doi:10.1128/JB.184.21.6026-6036.2002 (2002).
- 637 72 Goh, S., Chang, B. J. & Riley, T. V. Effect of phage infection on toxin production by  
638 *Clostridium difficile*. *J Med Microbiol* **54**, 129-135, doi:10.1099/jmm.0.45821-0 (2005).
- 639 73 Miles-Jay, A. *et al.* Longitudinal genomic surveillance of carriage and transmission of  
640 *Clostridioides difficile* in an intensive care unit. *Nat Med*, doi:10.1038/s41591-023-  
641 02549-4 (2023).
- 642 74 Seekatz, A. M. *et al.* Presence of multiple *Clostridium difficile* strains at primary infection  
643 is associated with development of recurrent disease. *Anaerobe* **53**, 74-81,  
644 doi:10.1016/j.anaerobe.2018.05.017 (2018).
- 645
- 646
- 647

648 **Figure Titles and Captions**



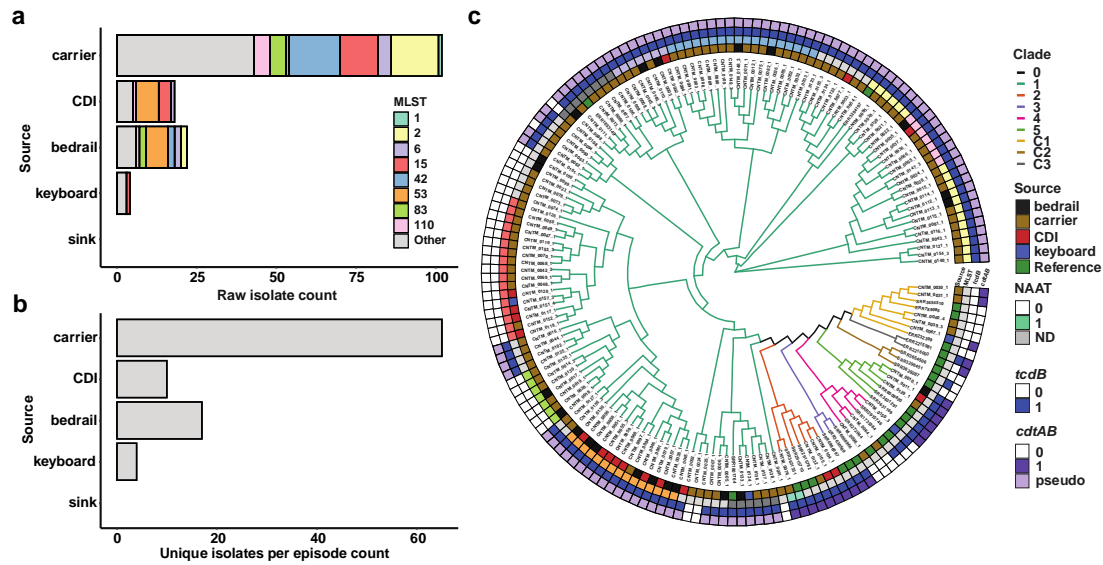
649

650 Figure 1: Study sampling and testing overview.

651 Caption: a) We sampled a leukemia and hematopoietic stem cell transplant ward at Barnes-  
652 Jewish Hospital in St. Louis, USA for 6 and 4 months respectively. Patients were enrolled and  
653 sampled upon admission, and then weekly for their time in the study wards. Surfaces were  
654 sampled weekly across the duration of the study. All samples and stool collected as part of  
655 routine clinical care were subjected to selective culture and MALDI-TOF MS identification, and  
656 isolates were whole-genome sequenced. Results of EIA testing as part of routine care were  
657 obtained.

658

659



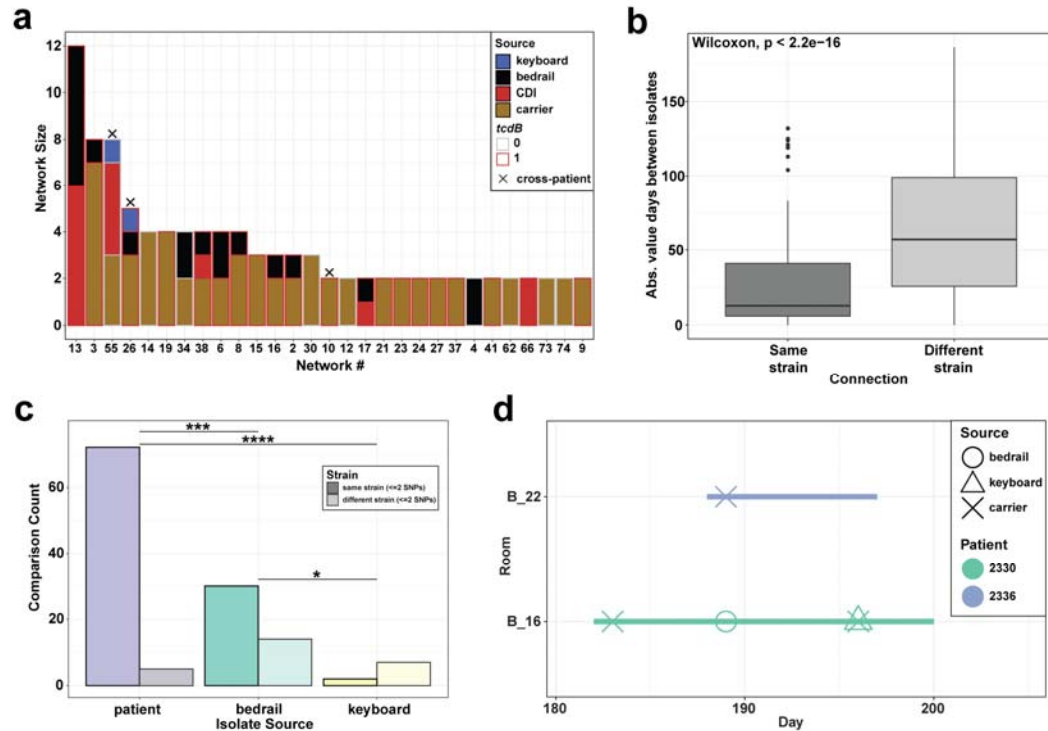
660

661 Figure 2: Total samples collected and phylogenetic relationships reveal carriers outnumber CDI  
662 patients and bedrails are the most commonly contaminated surface.

663 Caption: Total a) isolates collected and b) culture-positive episodes from each source. We found  
664 more carriers than CDI patients, and bedrails yielded the most *C. difficile* isolates. c) Cladogram  
665 of all isolates collected during this study plus references.

666

667



668

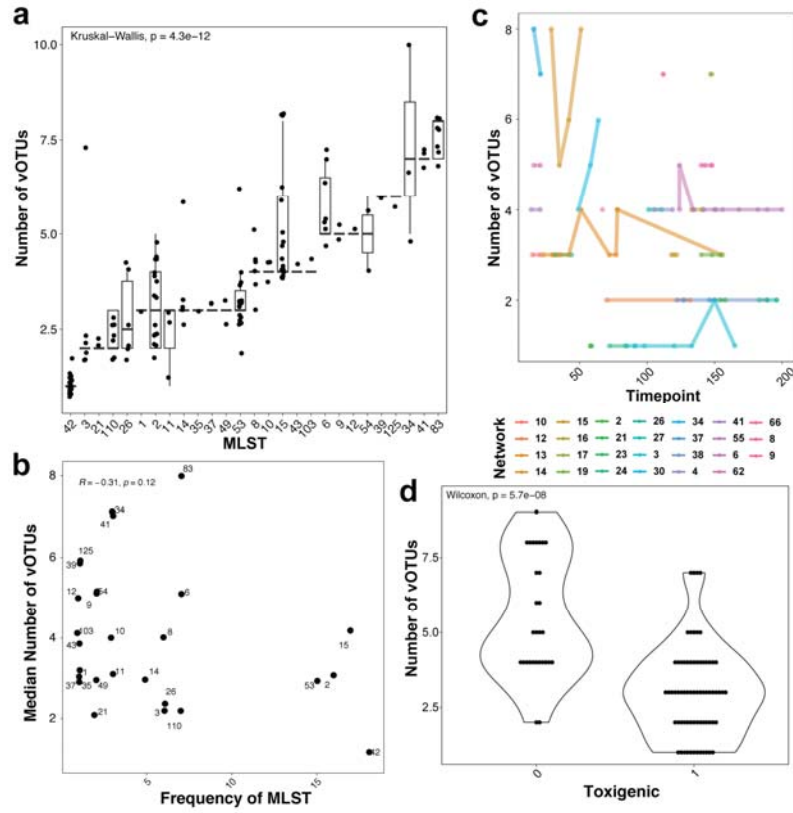
669 Figure 3: Hospital bedrails are a site of environmental contamination from colonized and CDI  
 670 patients.

671 Caption: a) Strain networks were defined by  $\leq 2$  core gene SNP cutoff. Network 55 includes the  
 672 non-toxicigenic isolates from Patient 2245 that are likely not responsible for the CDI. b) Absolute  
 673 value of days between isolates within strains and between strains. Isolates within the same  
 674 strain were significantly temporally linked ( $p < 2.2e-16$ , Wilcoxon test). c) Number of comparisons  
 675 in each group that fall within strain cutoff. Fisher's exact test, BH corrected. d) Strain tracking  
 676 diagram of transmission network 26, colors indicate patients and horizontal lines indicate stay in  
 677 a room. Patient 2336 sheds *C. difficile* onto the bedrail in room B\_16, and patient 2330 later is  
 678 identified as a carrier of the same strain.

679

680

681



682

683 Figure 4: Phage persistence in circulating *C. difficile* networks.

684 Caption: a) Phage diversity measured by phage population abundance for each isolate within an

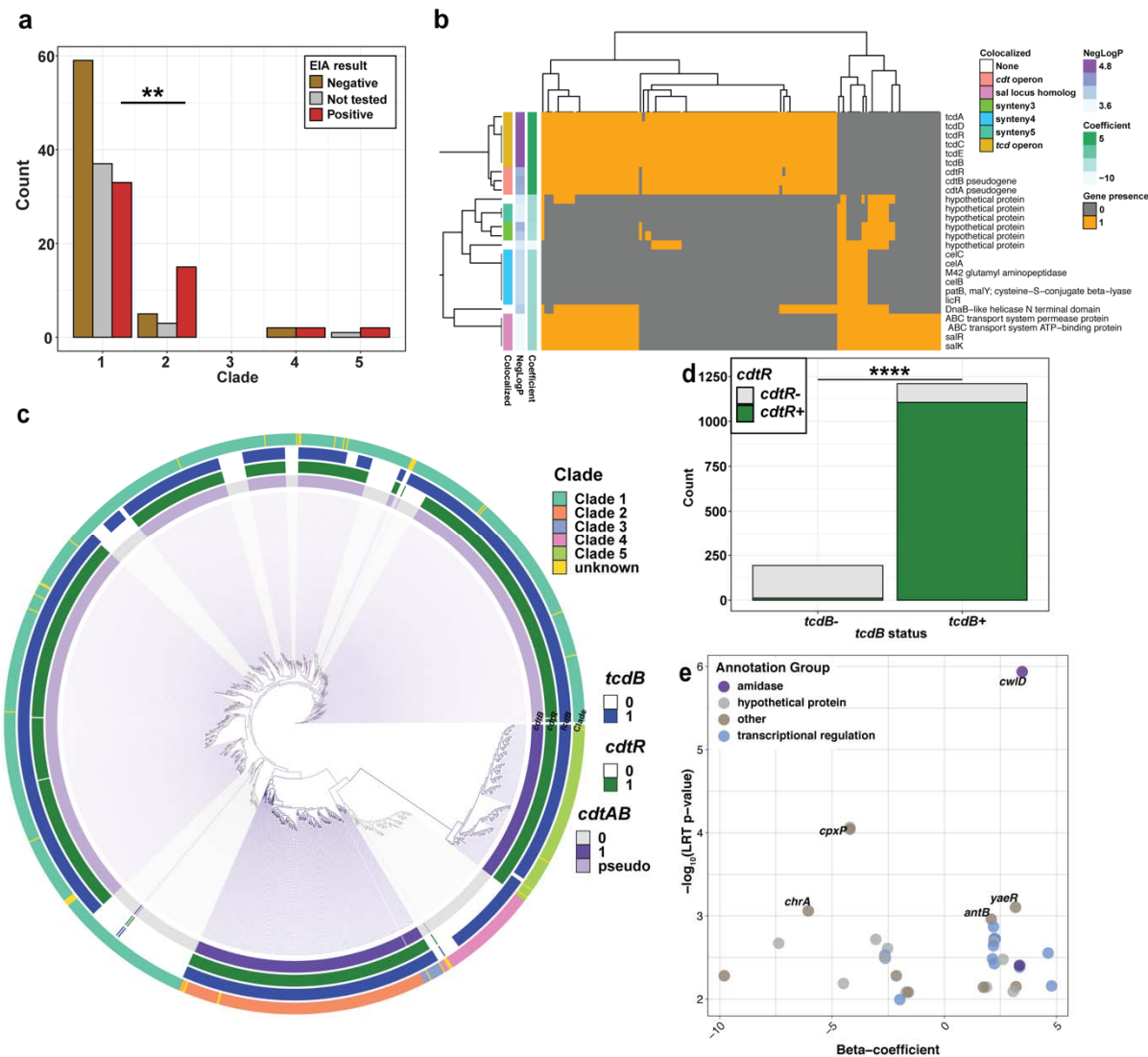
685 MLST. b) Relationship between phage diversity and frequency of ST in our cohort c) Temporal

686 trajectory of phage diversity for each network over time. d) phage population richness across

687 toxigenic and nontoxigenic isolates in our cohort, Wilcoxon test,  $p < 0.001$ .

688

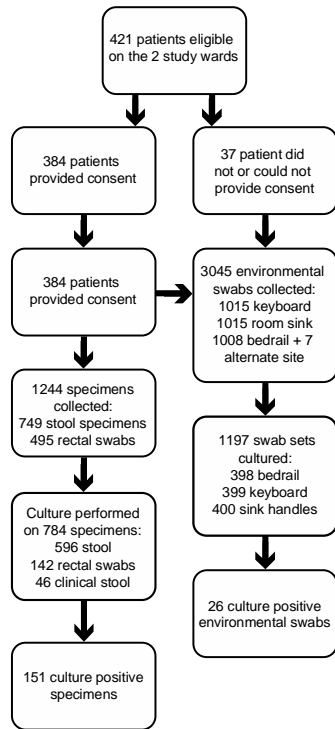




689  
 690 Figure 5: Clade 1 is responsible for the majority of CDI cases and carries unique correlates to  
 691 symptom severity.  
 692 Caption: a) EIA status by clade across this and a previous study<sup>29</sup>. Fisher's exact test,  $p < 0.01$ .  
 693 c) Phylogenetic tree of >1400 *C. difficile* isolates from NCBI (Supplementary Table 4) depicting  
 694 presence of binary toxin and PaLoc operons. d) Presence of full-length *cdtR* and association  
 695 with *tcdB* presence. e) Filtered results ( $p$ -values  $< 0.01$ ) pyseer analysis evaluating gene  
 696 association with CDI suspicion in Clade 1 isolates using the phylogenetically-corrected  $p$ -values  
 697 (LRT). Purple color indicates  $p < 0.001$ . Positive beta coefficient indicates gene association with  
 698 CDI suspicion, while negative beta indicates asymptomatic colonization.

699  
 700  
 701

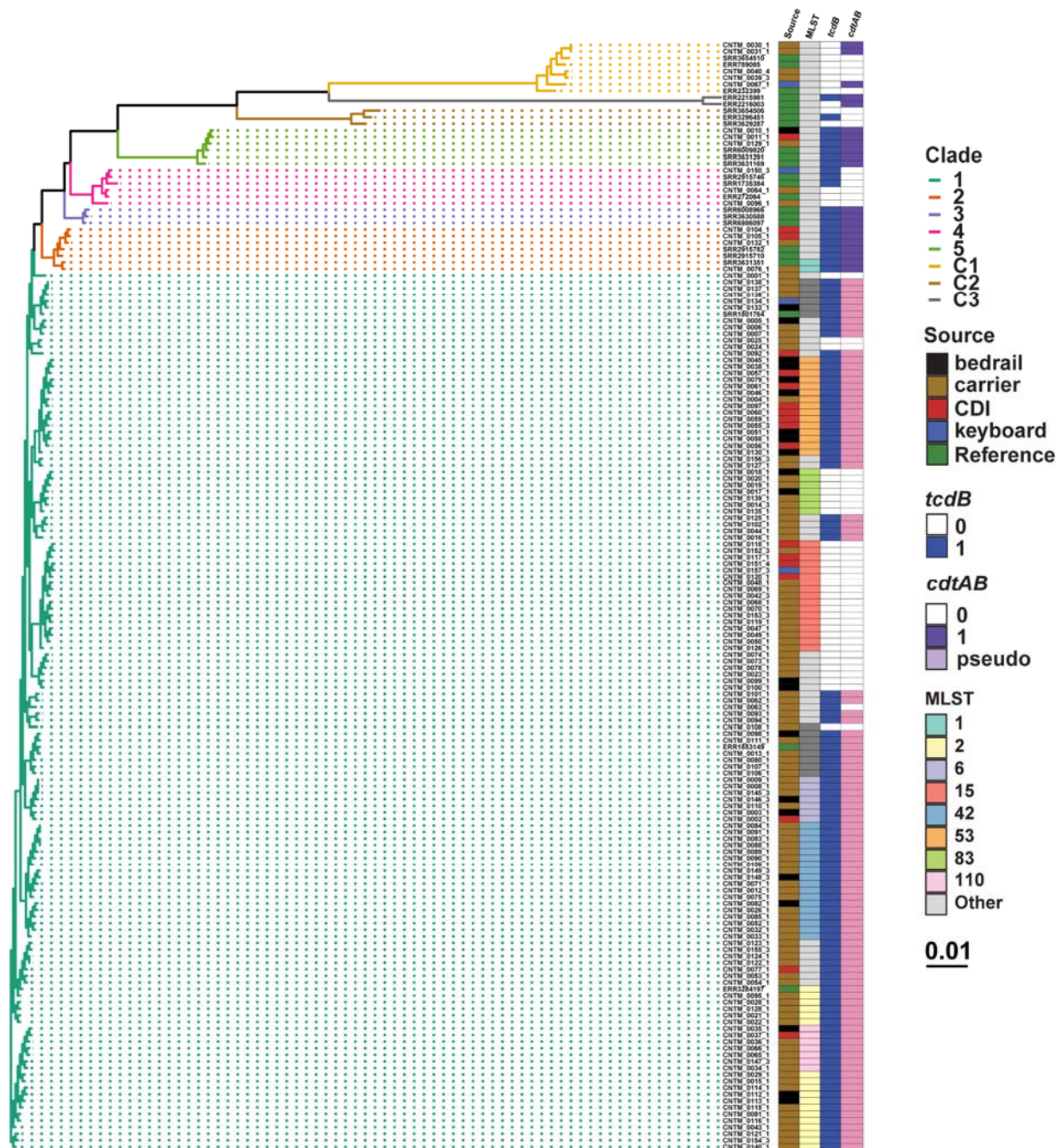
702 **Supplemental Figures and Tables Titles and Captions**

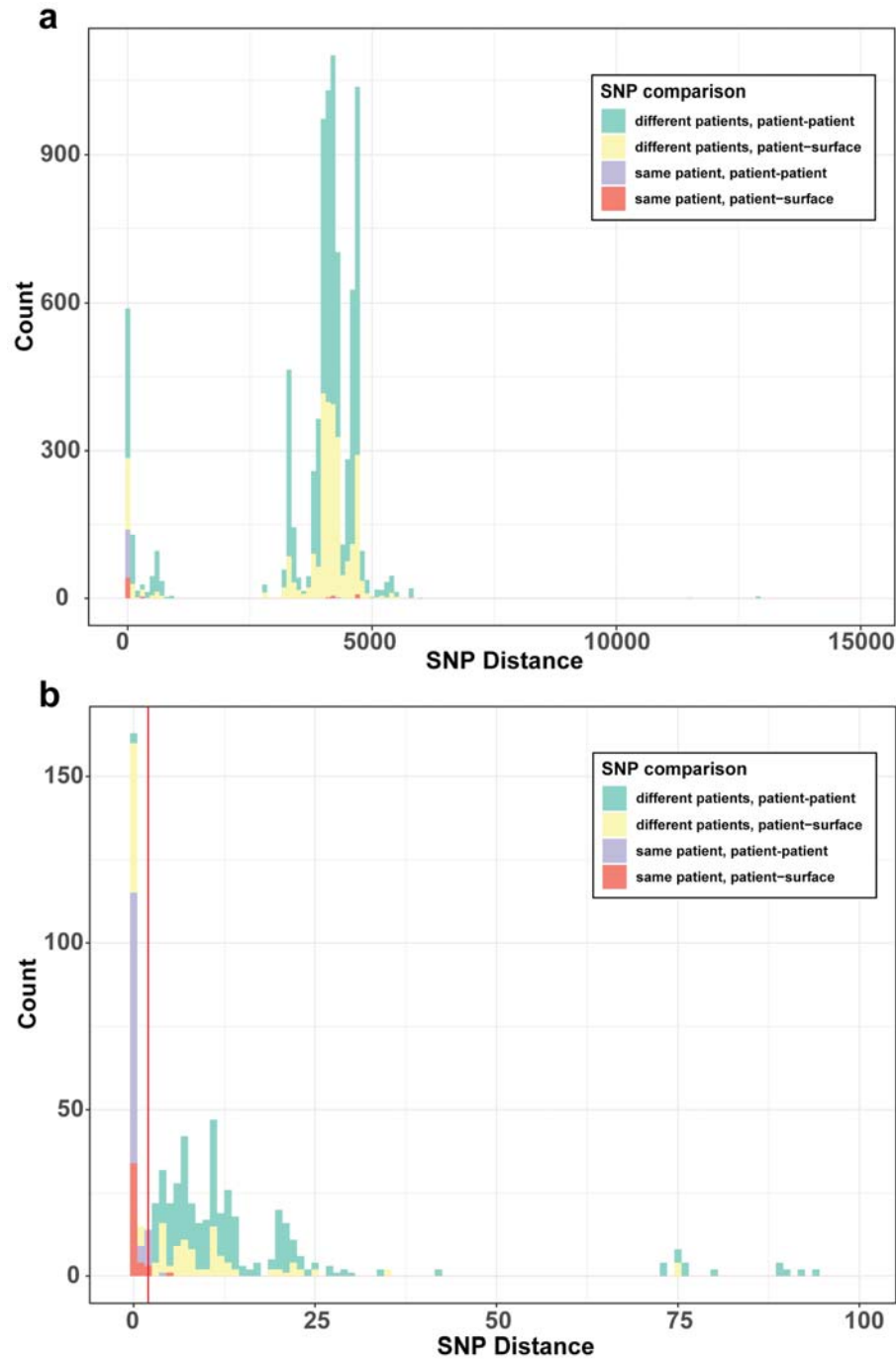


703

704 **Supplementary Figure 1: Bubble plot of enrollment, collection, and culture numbers.**

705  
706 Supplementary Figure 2: Phylogenetic tree of isolates collected in this study and select  
707 references (Supplementary Table 2).





708  
709 Supplementary Figure 3: Histogram of core genome SNP distances between different isolate  
710 comparisons, a) full histogram and b) zoomed to <200 SNPs.