

1 **PathoSPOT genomic surveillance reveals under the radar outbreaks of methicillin**
2 **resistant *S. aureus* bloodstream infections**

3 Ana Berbel Caban, MD ^{1*}, Theodore Pak, MD, PhD ^{2*}, Ajay Obla, PhD ², Amy Dupper, MA MPH
4 ¹, Kieran I. Chacko, PhD ², Lindsey Fox, MD ¹, Alexandra Mills, MD ¹, Brianne Ciferri, MPH²,
5 Irina Oussenko, PhD ², Colleen Beckford, MS ², Marilyn Chung, MS ², Robert Sebra, PhD ^{2,3,4,5},
6 Melissa Smith, PhD ^{2,3}, Sarah Conolly ⁶, Gopi Patel, MD ^{1,6}, Andrew Kasarskis, PhD ^{2,3,7},
7 Mitchell J. Sullivan, PhD ², Deena R. Altman, MD, MS ^{#1, 2}, Harm van Bakel, PhD ^{#2,3}

- 8 1. Department of Medicine, Division of Infectious Diseases, Icahn School of Medicine at
9 Mount Sinai, New York City, NY, USA
10 2. Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount
11 Sinai, New York City, NY, USA
12 3. Icahn Institute for Data Science and Genomic Technology, Icahn School of Medicine at
13 Mount Sinai, New York, NY 10029, USA
14 4. Black Family Stem Cell Institute, Icahn School of Medicine at Mount Sinai, New York,
15 NY, 10029, USA
16 5. Sema4, a Mount Sinai venture, Stamford CT, 06902
17 6. Infection Prevention, The Mount Sinai Hospital, New York City, NY USA
18 7. Department of Population Health Science and Policy, Icahn School of Medicine at Mount
19 Sinai, New York City, NY, USA

20 * These authors contributed equally to this work

21 # These authors are co-senior authors on this work

22 **Correspondence:**

23 Harm van Bakel (harm.vanbakel@mssm.edu)

24 **Keywords**

25 whole genome sequencing, visualization toolkits, MRSA bacteremia

26 Abstract

27 Whole genome sequencing (WGS) is increasingly used to map the spread and transmission
28 dynamics of human pathogens, especially in nosocomial settings. A limiting factor for more
29 widespread adoption of WGS for infection prevention practices is the availability of standardized
30 tools for genome analysis. Here we present PathoSPOT, which integrates genomic and
31 electronic medical record (EMR) data for rapid detection of nosocomial outbreaks and analysis
32 of epidemiological timelines. To demonstrate its capabilities we analyzed complete genomes
33 obtained from long-read sequencing of 197 invasive MRSA (bacteremia) cases from two
34 hospitals during a two-year period. PathoSPOT identified 8 clonal clusters that had not been
35 identified by epidemiological data, encompassing 33 patients (16.8%). The largest cluster was
36 consistent with an outbreak of a hospital-associated ST105 MRSA clone among 16 adult
37 patients in multiple wards. Timeline and location data suggested that an early outbreak event
38 led to infection and long-term colonization of multiple cases, followed by transmissions to other
39 patients during subsequent hospitalizations. Overall, we find that PathoSPOT can detect
40 outbreaks that are not readily apparent from epidemiological data and contribute to morbidity
41 and mortality. When applied as part of comprehensive ‘real-time’ sequencing-based
42 surveillance, PathoSPOT may play an important role in understanding the complexities of
43 pathogen transmission dynamics and enable timely interventions.

44 Introduction

45 The utility of whole genome sequencing to track transmissions and outbreak events is
46 well-established, in particular for highly clonal pathogens such as *S. aureus*, where classical
47 molecular typing methods such as multi-locus sequence typing (MLST) and pulsed-field gel
48 electrophoresis (PFGE) do not provide enough resolving power (Eyre et al. 2012; Harris et al.
49 2010; Köser et al. 2012; Sullivan et al. 2019). Despite the increased use of WGS, bottlenecks
50 remain that complicate its use in managing nosocomial outbreaks. Comparative genome
51 analyses often require specialized knowledge and/or selection of appropriate reference
52 sequences. Analysis and visualization frameworks are available to aid genome analysis in
53 outbreak settings (Hadfield et al. 2018; Zhou et al. 2018; Lees et al. 2019), but these are less
54 suited for nosocomial settings where genomic results need to be integrated with detailed
55 epidemiological results and patient histories to trace movements and contacts. This can be
56 time-consuming, especially when relying on manual chart review. Integration with electronic
57 medical record (EMR) systems can facilitate this process, but tools that combine
58 epidemiological and genomic information in a comprehensive manner for nosocomial outbreak
59 tracking are not readily available.

60 To facilitate detection and visualization of transmissions in nosocomial settings we developed an
61 open-source toolkit called PathoSPOT (Pathogen Sequencing Phylogenomic Outbreak Toolkit),
62 which combines automated comparisons of complete or draft genomes with interactive
63 visualization of clonal clusters. Further integration of epidemiological data enables high
64 resolution analysis of outbreak phylogenies and timelines. Since methicillin resistant
65 *Staphylococcus aureus* (MRSA) is a common cause of healthcare-associated (HA) infections in
66 the USA that pose a potentially fatal threat to patients, we used our toolkit as part of a complete
67 genome screening program of MRSA blood isolates, based on long-read sequencing. During a
68 24-month period, we obtained finished-quality MRSA genomes from 197 bacteremic patients at
69 two urban hospitals in New York City. PathoSPOT analysis identified 8 clusters of highly related
70 isolates that had not been recognized before based on available epidemiological information;
71 the most extensive cluster encompassing 16 patients. In-depth analysis with PathoSPOT
72 allowed us to reconstruct the outbreak timeline and identify common links among these
73 individuals. Our findings demonstrate the utility of PathoSPOT for infection prevention efforts in
74 health systems and highlight the role of colonization in long-term nosocomial outbreaks.

75 Results

76 MRSA surveillance with PathoSPOT identifies frequent under the radar transmissions

77 We developed PathoSPOT to enable automatic comparisons of large numbers of complete or
78 draft microbial genomes, and to rapidly identify closely related isolates indicative of transmission
79 events and map their epidemiological timelines (**Fig. 1A**). PathoSPOT combines existing tools
80 for whole genome comparison and alignment with custom analysis and visualization code
81 developed in Ruby, Python, and Javascript (<https://pathospot.org>). To demonstrate its utility, we
82 applied PathoSPOT to whole-genome surveillance data obtained from all MRSA isolates from
83 bacteremic cases at two hospitals (A and B) during a 2-year period. In total we sequenced 224
84 genomes for 221 isolates from 197 patients using PacBio long-read technology and obtained
85 184 finished-quality and 40 draft genome sequences (**Table S1**). In most cases we only
86 sequenced the primary blood culture, but additional isolates were sequenced for the same
87 patient in cases of prolonged or recurrent infections. We first used the PathoSPOT “compare”
88 pipeline to cluster genomes based on their Mash distance (Ondov et al. 2016). This step groups
89 highly related genomes prior to multi-genome alignment and avoids the need for manual
90 selection of a reference genome. The Mash distance threshold for MRSA was determined
91 empirically to yield clusters of genomes consistent with clonal complex assignments based on
92 MLST data, and to maximize core genome alignments in each cluster (**Fig. 1B**). Pairwise
93 distances between genomes were then calculated as the number of single nucleotide variants
94 (SNVs) between core genome alignments in each cluster for further analysis.

95 To identify transmission events we used the PathoSPOT “heatmap” tool (**Fig. 2**). We set a
96 threshold of ≤ 15 SNVs to identify potential transmissions, based on the extent of intra- and
97 inter-patient variability we previously observed in complete genome analysis of an extended
98 outbreak (Sullivan et al. 2019), and considering a core genome mutation rate of ~ 3 SNVs per
99 Mb per year (Harris et al. 2010; Young et al. 2012). The distance threshold can be varied
100 interactively in the heatmap tool to explore grouping at different levels of relatedness, depending
101 on the pathogen. At the selected threshold we identified 8 clonal clusters with a total of 33
102 patients, amounting to an overall transmission rate of 16.8% (**Fig. 2C**). Most clusters consisted
103 of patient pairs (5/8) but there were 3 with more than two patients. Patients within each cluster
104 typically had overlapping visits (75%) and stayed in the same wards at some point during these
105 visits (63%), but in many cases MRSA bacteremia was only found after they had transferred to

106 different wards. This likely contributed to the fact that none of the clusters had been recognized
107 epidemiologically. The most striking example of this was a clonal cluster of 16 patients, based
108 on data from 24 isolates obtained from 9 wards over a period of 21 months.

109 PathoSPOT timeline highlights the role of colonization in prolonged MRSA outbreaks

110 The presence of a clonal MRSA cluster among 16 bacteremia cases was consistent with a
111 prolonged “under the radar” outbreak. The outbreak strain was a clone of the
112 hospital-associated USA100 lineage (*spa* type t002, MLST 105, staphylococcal cassette
113 chromosome *mec* type II) and resistant to fluoroquinolones, oxacillin, clindamycin, erythromycin
114 and gentamicin based on automatic microbroth dilution testing. We next used the PathoSPOT
115 “dendro-timeline” tool, which combines phylogenetic analysis of outbreak isolates with the full
116 admission/transfer/discharge (ADT) history for each patient, to analyze this outbreak in more
117 detail (**Fig. 3**). The PathoSPOT core genome dendrogram, derived from the multigenome
118 alignment of the superset of isolates in the same Mash cluster, indicated the presence of distinct
119 subclades within the outbreak (**Fig. 3A**). Shared variant profiles within each subclade were
120 consistent with sub-transmissions within the larger outbreak (**Fig. 3A**). Isolates from p176, who
121 tested positive for MRSA bacteremia numerous times within a span of 6 months, were
122 represented in distinct clades, suggesting that this patient carried multiple variants of the
123 outbreak strain at the same time. This was confirmed by sequencing two subclones from p176
124 isolate ER05682 (Fig. 3A, ▲ and ◆).

125 Based on the PathoSPOT timeline of events (**Fig. 3B**), we reconstructed the most likely
126 outbreak scenario (**Fig. 3C**). The first two cases (p40 and p574) tested positive on ward 8. No
127 other positives were found on this ward, but five other cases had overlapping stays (p176, p459,
128 and p181) or were admitted to the same ward within four weeks of the first positive test (p34,
129 p142). All but one patient tested positive for bacteremia within seven weeks of their stay in ward
130 8, following transfers or readmissions to other wards. Strikingly, p459, who was discharged from
131 ward 8 four days after the initial positive case, did not present with bacteremia until 20 months
132 later. In the intervening period this patient had no contact with our health system except for two
133 dermatology office visits where positive wound cultures for MRSA were obtained. The degree of
134 genetic drift of the p459 isolate (11 SNVs) and pattern of positive wound cultures prior to
135 bacteremia are consistent with long-term colonization after initial exposure in ward 8, although
136 we could not verify this scenario as the wound isolates were not available for sequencing.

137 Five additional patients tested positive in wards 16, 2, and 12 during the first six months of the
138 outbreak (**Fig. 3B, 3C**). Each instance was preceded by the transfer of an outbreak case that
139 had previously stayed in ward 8, suggesting that direct or indirect transmissions from these
140 cases propagated the outbreak. Patient p34 was transferred from ward 8 to ward 16, into a
141 room neighboring p628, who became bacteremic two days later. Both their isolates were
142 grouped in the same subclade (**Fig. 3A, ▼**). Likewise, p142 was transferred from ward 8 to
143 ward 2, where there was overlap with p399 and p476 before all three became bacteremic on
144 this ward. Notably, p142 and p176 overlapped with p476 on two different days in the inpatient
145 hemodialysis unit, providing an alternative acquisition route. Finally, after discharge from ward 8,
146 p181 was readmitted to ward 12, where the patient overlapped stays with patients p2 and p669.

147 From months 7 to 21, four additional transmissions were identified. Patient p176 tested positive
148 for the outbreak strain on multiple occasions during readmissions and was implicated in two
149 other potential transmissions. The patient visited the emergency department (ED) on the same
150 day as patient p648, and had an overlapping stay in ward 2 with p593 for at least 5 days, in the
151 months prior to their positive tests. Following readmission after a 20-month hiatus, p459 likely
152 transmitted to p10, based on evidence of an overlapping stay in ward 36 (**Fig. 3B**) and the high
153 relatedness of their isolates (**Fig 3A, +**). Patient p77 is the only person that did not have
154 overlapping stays with other outbreak cases. The patient had a total of two pediatric (ward 49)
155 and one adult (ward 2) admission to the hospital within 21-months. Given that all other outbreak
156 cases were adults, we consider ward 2 the most likely location of MRSA acquisition, where p77
157 shared healthcare workers with p593 who was admitted to the same unit 11 weeks before.

158 Altogether, PathoSPOT analysis suggested that initial exposure in Ward 8 resulted in
159 colonization and clinical infection of up to 7 cases (44% of the prolonged outbreak cluster),
160 followed by secondary transmissions after ward transfers and/or subsequent hospital visits of
161 these initial cases. The alternative scenario of community transmissions was assessed by
162 mapping of home zip codes, which indicated that 13 of 16 cases lived in different areas,
163 including 5 that lived outside of New York City. Spatiotemporal analyses of the seven smaller
164 clusters (**Fig. S1**) showed that five include direct overlaps in wards or outpatient settings, of
165 which three are plausible transmission events, and two such events are months before the
166 clonal blood cultures were obtained. Taken together, our data suggest that in the absence of
167 genomic surveillance, MRSA transmission within a healthcare setting often goes undetected as
168 asymptomatic colonization, only leading to clinically apparent infections weeks to months later.

169 Hand hygiene compliance and vascular access implicated in under-the-radar outbreak

170 As the prolonged outbreak extended over multiple wards, we further investigated hand hygiene
171 rates, shared HCW, patient movements, and clinical characteristics. Average hand hygiene
172 compliance in the wards associated with the majority of cases ranged between 79-83%.
173 Compliance in wards 8 and 16 was reduced to 70% and 66% in the month prior to the first
174 outbreak case, while ward 2 compliance was maintained at 79%. All outbreak patients shared at
175 least one HCW involved in the care of other patients in the cluster. This is consistent with the
176 high degree of overlapping stays in the same ward and suggests that direct and indirect
177 transmissions may have played a role in propagating the outbreak. Although outbreak cases
178 were moved frequently between units based on transfer records, they did not move more
179 frequently than non-outbreak patients.

180 Chart review of outbreak cases revealed that 69% (n=11) were male, 75% (n=12) had been
181 admitted from home, and 63% (n=10) had a hospital admission in the 90 days prior (**Table 1**).
182 75% (n=12) were considered hospital-onset (HO)-MRSA as defined by NHSN (Center for
183 Disease Control and Prevention 2018) and 88% (n=14) of subjects had an invasive device at
184 the time of infection. Notable was the presence of active malignancy (57%; n=8) including five
185 persons with leukemia (4 acute myeloid leukemia, 1 chronic lymphocytic leukemia), and others
186 with multiple myeloma (n=1), disseminated Kaposi sarcoma (n=1), and metastatic breast cancer
187 (n=1). Among patients with hematologic malignancies, three had undergone hematopoietic stem
188 cell transplant. The most common presumed source of bacteremia was vascular access (n=9;
189 56%), followed by skin source (n=4; 25%). The 90-day mortality incidence was 25% (n=4), of
190 which 75% (n=3) was related to bacteremia with outbreak strain. We performed univariate and
191 multivariate analyses of the 16 outbreak cases compared with 34 unique patients infected with
192 non-outbreak ST105 MRSA during the same time period. Outbreak patients were significantly
193 associated with HO-MRSA (OR= 5.20 95% CI [1.04-26.01]; p=0.04) and the administration of
194 intravenous chemotherapy prior to bacteremia (OR=11.24 95% CI [1.72-73.28]; p=0.01) in
195 multivariate analysis. There were no differences in outcomes between outbreak and
196 non-outbreak patients.

197 Discussion

198 We developed the PathoSPOT phylogenomics toolkit to facilitate detection of outbreaks and
199 transmissions as part of an ongoing genomics-based pathogen surveillance program.
200 Application of the toolkit to 197 patients with MRSA bacteremia over a two-year period shows
201 that under-the-radar transmission events are important sources of morbidity and mortality. Since
202 precise construction of epidemiological timelines is labor intensive and time consuming, when
203 implemented, our toolkit can avoid these bottlenecks and assist infection prevention staff with
204 rapid identification of transmission events.

205 During the study period, PathoSPOT detected 8 clusters of MRSA involving 33 patients with
206 bacteremia. A large outbreak impacted 16 patients from distinct adult medicine wards and
207 spanned nearly the entire study period. In the absence of routine patient, HCW and
208 environmental screening it was not possible to determine the original source of transmission
209 cascade. Ward overlap was a key factor in forward transmissions, with a delay between
210 exposure and subsequent bacteremia that sometimes extended for several months and likely
211 played a role in obscuring the outbreak from epidemiological detection. Most wards involved
212 had lower hand hygiene rates surrounding key outbreak events. Additional contributing factors
213 likely included frequent patient movements and shared HCW. Frequent room changes within
214 and between wards may also have resulted in contaminated environmental surfaces, which has
215 been shown to play a role in the transmission of hospital pathogens (Otter et al. 2011, 2013).

216 Our study suggests a role for colonization in the progression and persistence of long term
217 outbreaks (Senn et al. 2016; Merrer et al. 2000). Skin colonization may have played a role in
218 acquisition, as the majority of presumed patient sources of bacteremia were central venous
219 access or skin source. The number of patients with malignancies was also notable, particularly
220 hematological with bone marrow suppression. Patients with hematologic malignancies have an
221 increased risk of bacteremia, as central venous catheters remain an essential tool for their
222 treatment, frequently leading to catheter-related infections (Zakhour et al. 2016).

223 The detection of nosocomial transmissions and outbreaks is critical for healthcare organizations.
224 Precise, genomics-based, pathogen surveillance programs supported by rapid analysis
225 frameworks such as provided by PathoSPOT are essential to detect events that are not readily
226 ascertained by conventional infection prevention methods. Such programs would be more

227 effective when implemented across regional health systems, long-term acute care hospitals, and
228 skilled nursing facilities, to track dissemination of strains and identify sources and at-risk
229 patients based on contact networks. When combined with timely intervention, these efforts could
230 help reduce endemic rates of nosocomial infections.

231 Methods

232 Ethics statement

233 This study was reviewed and approved by the Institutional Review Board of the Icahn School of
234 Medicine at Mount Sinai.

235 Bacterial culturing, DNA extraction and sequencing

236 Isolates were cultured and identified as part of standard clinical testing procedures in the Mount
237 Sinai Hospital Clinical Microbiology Laboratory (CML), and stored in tryptic soy broth (TSB) with
238 15% glycerol at -80°C. Selected isolates were subcultured on tryptic soy agar (TSA) plates with
239 5% sheep blood (blood agar) (ThermoFisher Scientific) under nonselective conditions. The
240 Qiagen DNeasy Blood & Tissue Kit (Qiagen) was used for DNA extraction, as previously
241 described (Sullivan et al. 2019). Following quality control, DNA and library preparation,
242 long-read sequencing was performed on the Pacific Biosciences RS-II platform to a depth of
243 >200x.

244 Genome assembly

245 PacBio SMRT sequencing data were assembled using a custom genome assembly and
246 finishing pipeline (<https://github.com/powerpak/pathogendb-pipeline>), as previously described
247 (Sullivan et al. 2019).

248 Comparative genome analysis using PathSPOT-compare

249 We developed the PathoSPOT-compare pipeline to perform comparative phylogenomic analysis
250 of finished, annotated genome assemblies for the specific purpose of outbreak detection. The
251 pipeline is implemented as a Rakefile (a Makefile for the Ruby language) that calculates

252 dependencies and executes all necessary subtasks to reach desired outputs.
253 PathoSPOT-compare takes FASTA-formatted genome assemblies as input, along with a
254 relational database (SQLite or MySQL) containing metadata for each assembly (including
255 collection time, location, collection method, organism, and patient ID), as well as metadata on
256 patient ADT history (for spatiotemporal analysis).

257 Genetic distances for outbreak detection are ultimately calculated by counting SNP differences
258 within core genome alignments; however, there is a trade-off between aligning increasingly
259 diverse assemblies and a diminishing core genome size (as more subsequences will fail to align
260 across all assemblies). Therefore, we implemented a hybrid approach, wherein pairwise
261 distances between all assemblies are first estimated using Mash (Ondov et al. 2016), which
262 uses a k-mer based hashing approach that approximates average nucleotide identity (ANI).
263 Mash distances are used to perform greedy single-linkage hierarchical clustering, with clusters
264 capped at a pre-specified diameter and size. The default parameters, which are also the
265 parameters used for this study, are a maximum Mash cluster diameter of 0.02 (approximating
266 98% ANI among all included genomes) and at most 100 genomes per cluster.

267 Rapid core genome alignments are then created for each cluster using parsnip (Treangen et al.
268 2014), which is tailored for intraspecific genome analysis and is therefore well-suited for
269 outbreak analysis. Outputted variant call files (VCF) for each cluster are converted to NumPy
270 arrays (NPZ files) for fast loading and subsetting of variant data by PathoSPOT-visualize, the
271 downstream visualization web application that can display called variants alongside
272 phylogenies. The primary output for PathoSPOT-visualize is a JSON file containing a matrix of
273 pairwise SNP distances for all genomes (with inter-cluster distances left unspecified) and a
274 maximum-likelihood phylogeny for each cluster. Additional optional pipeline tasks export patient
275 location data (as TSV files) and epidemiologic data on positive and negative culture results (as
276 JSON files), both of which are automatically utilized and layered onto the comparative genomic
277 analyses within PathoSPOT-visualize when available.

278 Interactive detection and visualization of outbreaks with PathoSPOT-visualize

279 To visualize the above analyses as depicted in **Fig. 2, 3A-B, S1, and S2**, we created the
280 PathoSPOT-visualize interactive web application. The application uses PHP scripts and AJAX to
281 serve data from the JSON, TSV, and NPZ output files generated by the PathoSPOT-compare
282 pipeline, which are then dynamically mapped to interactive HTML5 and scalable vector graphics

283 (SVG) elements using the D3.js (Data-Driven Documents) framework. All views are rendered in
284 the browser, allowing the user to alter settings that trigger live animated transitions, and an
285 intuitive sense of how changes propagate between the linked views of data.

286 There are three main user interfaces, the “heatmap” tool, the “network map” tool and the
287 “dendro-timeline” tool. Users initially interact with the “heatmap” tool, which starts with selection
288 of a dataset that can be prefiltered by specimen location, MLST, and time interval. The user can
289 dynamically adjust the SNP threshold that specifies the genetic distance deemed suspicious for
290 transmission; this process is aided by a histogram depicting distributions of SNP distances
291 among same-patient isolates (generally expected to be related) and different-patient isolates
292 (which are not, assuming a low level of transmission). This threshold is used to perform
293 single-linkage hierarchical clustering of genomes on the client-side, with the clusters assigned
294 random colors and depicted on a beeswarm timeline plot and a large heatmap of pairwise
295 distances among all selected genomes. This large heatmap can be swapped for the “network
296 map” view (**Fig. S2**), which plots genomes by their collection location in a geospatial layout,
297 overlaid with density plots of overall epidemiological incidence and force-directed network links
298 depicting genetic relationships.

299 Epidemiological links within the clusters can be further explored in the “dendro-timeline” tool,
300 which combines a traditional phylogenetic dendrogram with a SNP matrix, a mapping of SNP
301 locations onto a reference assembly, and a pannable-zoomable timeline of patient locations
302 over time, with spatiotemporal overlaps highlighted as bright arcs. The phylogeny for the
303 “dendro-timeline” tool is extracted from the larger maximum-likelihood trees built by parsnp,
304 based on the SNP threshold and clustering parameters that the user selected in the “heatmap”
305 tool.

306 Case review

307 We performed a retrospective clinical chart review on all adults (age ≥ 18) subjects identified
308 with MRSA bacteremia. Analyses were performed in SAS (v9.4)(SAS Institute 2013). Variables
309 were analyzed initially in a univariate logistic regression model. Variables $p \leq 0.2$ were then
310 placed into a stepwise multivariate logistic regression model (Dupper et al. 2019). An additional
311 in-depth chart review was performed for subjects identified as being involved in transmission
312 events. These details included location (ward, room, bed), all ADT information, procedures, and
313 providers. Whole genome sequencing was performed on 1st patient blood isolate positive for

314 MRSA as part of an ongoing genomic surveillance program as described in (Sullivan et al.
315 2019). Hand hygiene is monitored by the Infection Prevention and Control department by the
316 use of anonymous observers using the Joint Commision's Targeted Solutions Tool (TST)
317 (Anderson et al. 2019), which was implemented in November 2014. This tool allows the staff
318 member to document reasons for non-compliance and target areas of interventions. Hand
319 hygiene observations are collected anonymously at entry and exit by trained staff members in
320 each hospital ward. Hand hygiene rates were compared using a one-sample t-test.

321 Data availability

322 All study data is available at <https://www.pathospot.org>.

323 Software

324 The PathoSPOT-compare and PathoSPOT-visualize packages developed for this study are both
325 open source and can be obtained from:

326 - <https://github.com/powerpak/pathospot-compare>

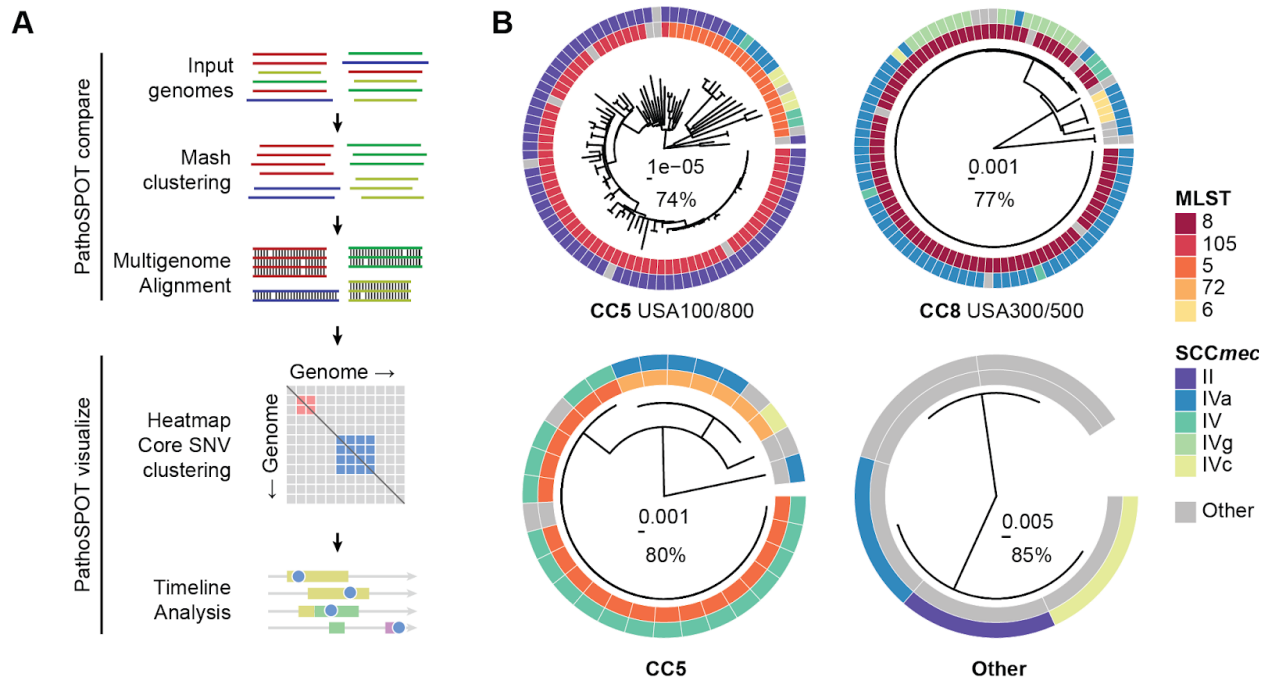
327 - <https://github.com/powerpak/pathospot-visualize>

328 A live demo of all visualizations created for this study, along with documentation on setting up
329 and using the software with example data from this study, can be found at <https://pathospot.org>.

330 Acknowledgements

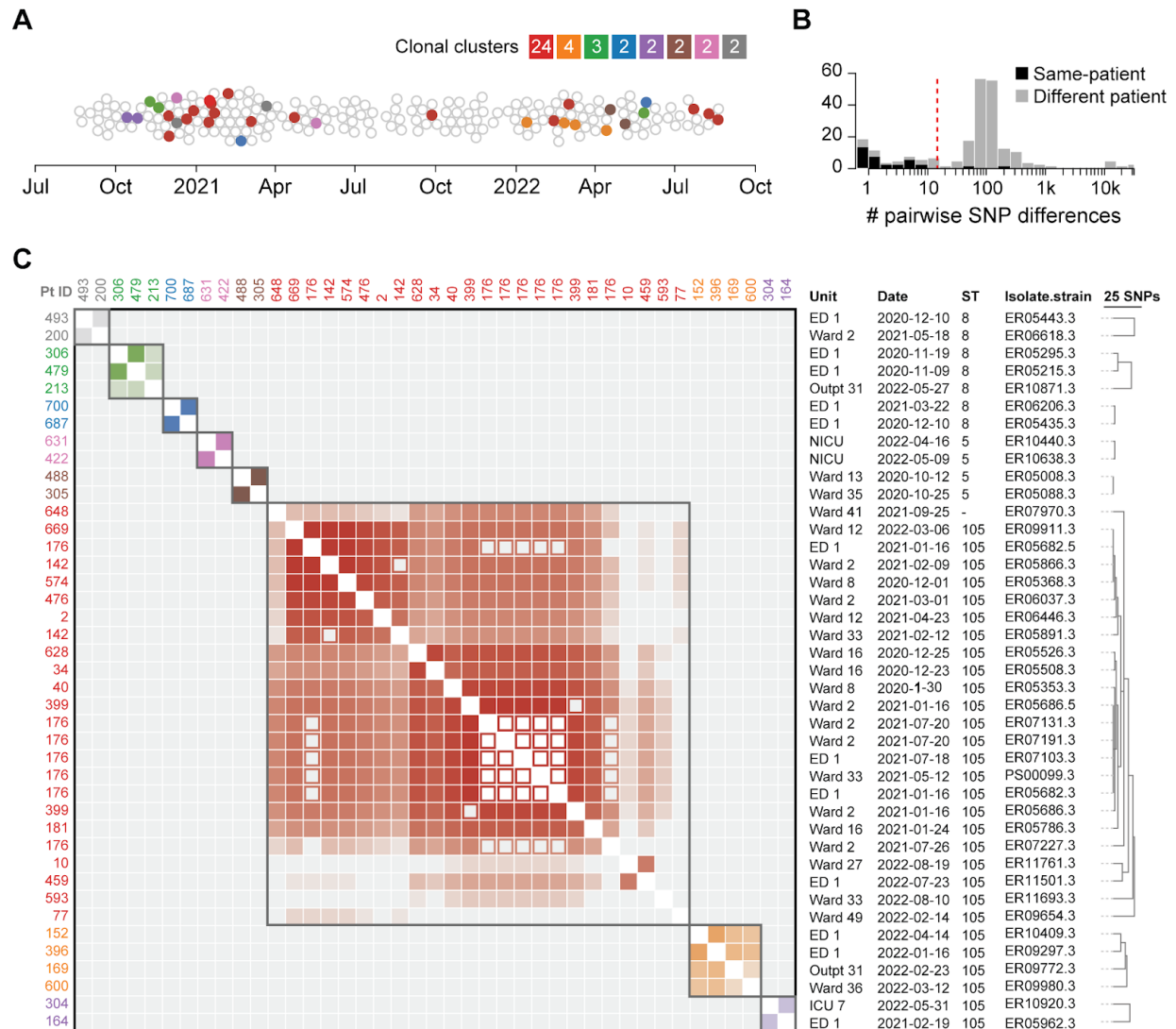
331 This research was supported in part by R01 AI119145 (H.v.B.), the Icahn Institute for Genomics
332 and Data Science (A.K.), the CTSA/NCATS KL2 Program (KL2TR001435, Icahn School of
333 Medicine at Mount Sinai; D.R.A), and the New York State Department of Health Empire Clinical
334 Research Investigator Program (Awarded to Judith A. Aberg, Icahn School of Medicine at Mount
335 Sinai; D.R.A.) and F30 AI122673 (T.R.P.). The Research reported in this paper was supported
336 by the Office of Research Infrastructure of the National Institutes of Health (NIH) under award
337 numbers S10OD018522 and S10OD026880 as well as institutional funds. The funders had no
338 role in study design, data collection and interpretation, or the decision to submit the work for
339 publication.

340 Figures



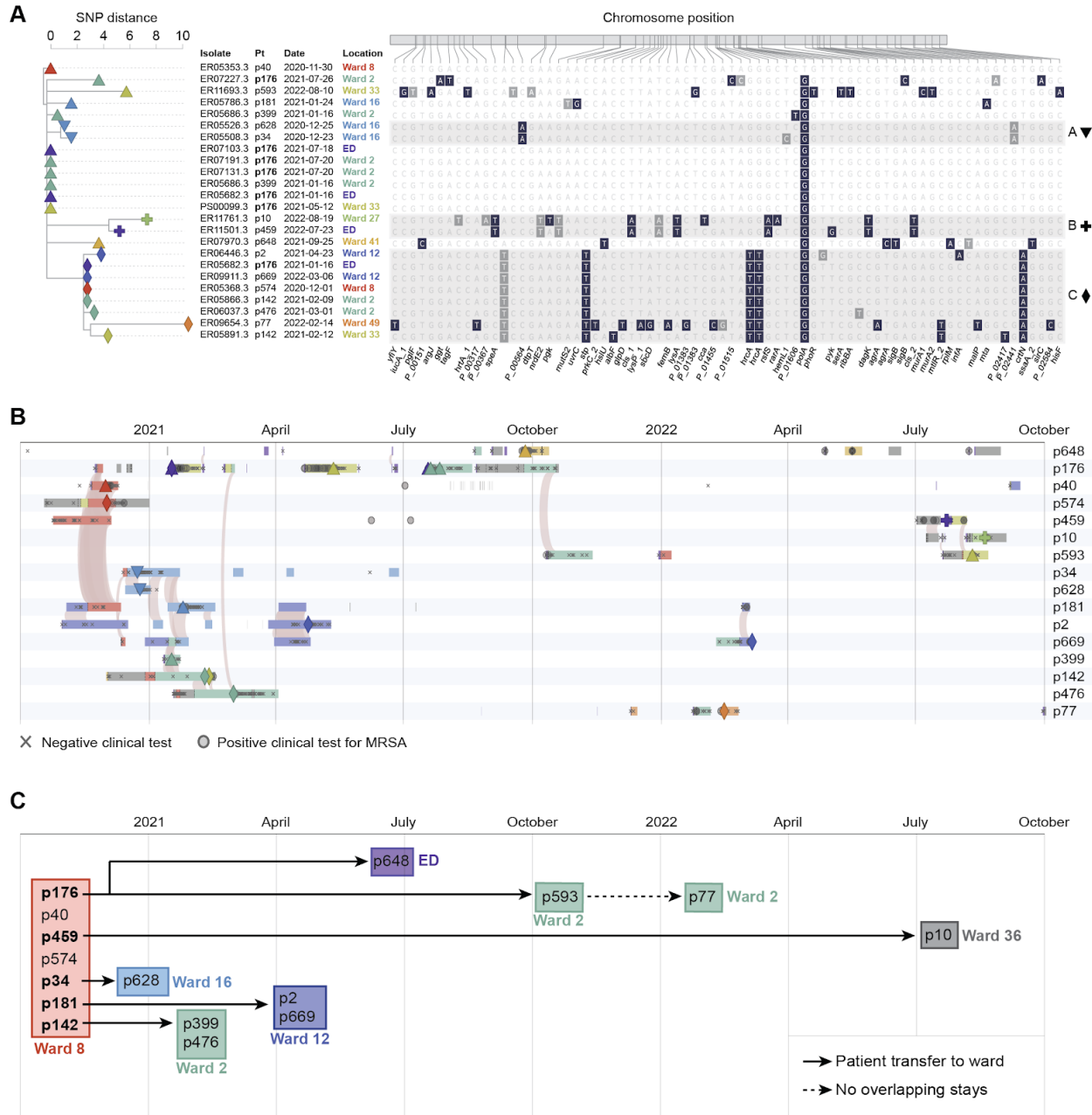
341 Figure 1. PathoSPOT comparative genome analysis of 221 MRSA isolates

342 **A)** Overview of the PathoSPOT whole-genome comparison framework. **B)** Maximum-likelihood
 343 phylogenetic trees produced from core genome SNVs identified from Parsnp whole-genome
 344 alignments of 4 clusters identified at a Mash score threshold of 0.02. Trees are annotated with
 345 MLST and *SCCmec* information (key shown on right) and clonal complexes (bottom). Scale
 346 bars indicating the number of substitutions per site in the phylogeny and the percentage of core
 347 genome coverage among all sequences is shown at the center of each tree.



348 Figure 2. Identification of clonal clusters among 197 MRSA bacteremia cases
 349 **A)** Beeswarm plot of MRSA cases with sequenced isolate genomes during the surveillance
 350 period. Cases with isolates separated by ≤ 15 core genome SNPs are grouped in clonal clusters,
 351 each highlighted with a distinct color. The number of isolates in each cluster is indicated in the
 352 color key. **B)** Histogram of pairwise core genome SNP distances for isolates obtained from the
 353 same patient (black bars) and isolates obtained from different patients (grey bars). The vertical
 354 red line indicates the 15-SNP threshold for clonality. **C)** Heatmap of pairwise core genome SNP
 355 distances between clustered isolates. Clusters are grouped along the diagonal and colored as in
 356 A, with decreasing shading reflecting an increased pairwise SNP distance. Closed squares and

357 open squares are used for isolates from different patients or the same patient, respectively. All
358 date information in this figure was recoded to protect health information.



359 Figure 3. Epidemiologic timeline of long-term outbreak.

360 **A**) Phylogenetic tree of core genome SNP differences (left) with corresponding locations of each
 361 variant relative to the first outbreak isolate that was obtained from p40 (right). Non-synonymous
 362 and synonymous variants are highlighted in black and grey, respectively. Distinct clades with
 363 three or more shared variants are shown as shaded areas (key on right). **B**) Epidemiologic
 364 timeline reconstructed by PathoSPOT integration of genomic and epidemiological data. Rows

365 correspond to patients, with admission periods in hospital wards shown as horizontal bars.
366 Wards are colored as in A. Sequenced clinical isolates are shown as different symbols matching
367 those used in the phylogenetic tree in A, and colored according to the collection ward. Shaded
368 arcs signify ward-level patient overlap within 24 hours. Other positive or negative clinical test
369 results are indicated by grey symbols, with a key shown below. The timeline scale is shown at
370 the top. **C)** Summary of key outbreak events in wards, derived from the epidemiological timeline.
371 See main text for details. All date information in this figure was recoded to protect health
372 information.

373 Table 1. Outbreak patients vs. patients with MLST 105 between January 2016 through
374 December 2017

Factor	Outbreak Patients N = 16 (%)	Non-Outbreak Patients N = 34 (%)	Univariate Analysis		Multivariate Analysis	
			OR (95% CI)	p value	OR (95% CI)	p value
Male	11 (69)	21 (62)	1.36 (0.39-4.82)	0.63		
<i>Race/Ethnicity</i>						
Non-Hispanic White	5 (31)	14 (41)	Reference			
Non-Hispanic Black	3 (19)	9 (26)	0.93 (0.18-4.90)	0.94		
Hispanic/Latino/Asian	5 (31)	6 (18)	2.33 (0.49-11.17)	0.29		
Unknown	3 (19)	5 (15)	1.68 (0.29-9.75)	0.56		
<i>Age at Time of Infection</i>						
18-54 Years	6 (38)	7 (21)	Reference			
55-69 Years	5 (31)	9 (26)	0.65 (0.14-3.04)	0.58		
≥ 70 Years	5 (31)	18 (53)	0.32 (0.07-1.41)	0.13		
<i>History of IV Drug Use</i>						
HIV	2 (13)	2 (6)	2.29 (0.29-17.90)	0.43		
<i>Admission Source</i>						
Home	1 (6)	3 (9)	0.69 (0.07-7.19)	0.76		
Home	12 (75)	17 (50)	Reference			
NH/Rehab/LTACH	2 (13)	11 (32)	0.26 (0.05-1.38)	0.11		
Other Hospital	2 (13)	6 (18)	0.47 (0.08-2.75)	0.40		
Prior Hospital Admission (90 Days)	10 (63)	23 (68)	0.80 (0.23-2.76)	0.72		
<i>NHSN Definitions</i>						
CO-MRSA	4 (25)	24 (71)	Reference		Reference	
HO-MRSA	12 (75)	10 (29)	7.20 (1.87-27.79)	0.004	5.20 (1.04-26.01)	0.04
^A Presence of Invasive Device	14 (88)	27 (90)	0.78 (0.12-5.21)	0.80		
^B Receiving Cancer Treatment	7 (44)	2 (6)	15.75 (1.75-141.39)	0.01	11.24 (1.72-73.28)	0.01
<i>Charlson Comorbidity Index (CCI)</i>						
0-3	4 (25)	9 (26)	Reference			
≥ 4	12 (75)	25 (74)	1.08 (0.28-4.23)	0.91		
<i>History of MRSA Colonization</i>						
6 (38)	6 (38)	16 (47)	0.68 (0.20-2.28)	0.53		
<i>Presumed Source of MRSA BSI</i>						
^C Skin Source	4 (25)	7 (21)	1.29 (0.32-5.24)	0.73		
Pneumonia	1 (6)	6 (18)	0.31 (0.03-2.83)	0.30		
[*] , [†] Vascular access	9 (56)	9 (26)	3.57 (1.03-12.43)	0.05		
Other/Undetermined	2 (13)	12 (35)	0.26 (0.05-1.35)	0.11		
Persistent Bacteremia (≥ 5 Days)	2 (13)	9 (26)	0.40 (0.08-2.10)	0.28		
ICU Admission Prior to BSI	4 (25)	3 (9)	3.44 (0.67-17.73)	0.14		
Intubated Prior to MRSA BSI	4 (25)	2 (6)	5.33 (0.86-33.00)	0.07		

375 **Bold** = significant at ≤ 0.05

376 Abbreviations: BSI, bloodstream infection; HIV, human immunodeficiency virus; ICU, intensive care unit; IV, intravenous; NHSN, National Healthcare Safety
377 Network.

378 ^A Includes devices such as pacemaker, any vascular access, orthopedic hardware, foley catheter, arteriovenous graft placement, percutaneous endoscopic
379 gastrostomy (PEG), ostomy, or any type of urinary collection at the time of first positive bloodstream infection.

380 ^B Includes patients actively receiving cancer treatment through a central venous catheter prior to bacteremia in the outpatient or inpatient setting.

381 ^C Skin source includes skin and soft tissue infections, thrombophlebitis due to peripheral IV catheters

382 ^{*} Variable not included in the multivariate analysis in order to prevent collinearity between receiving cancer treatment and vascular access.

383 [†] Vascular access devices include: non-tunneled central venous catheter, tunneled catheter (hickman or permacath), implanted port, peripherally inserted central
384 catheter (PICC line), as well as arteriovenous graft (AVG) and fistula (AVF).

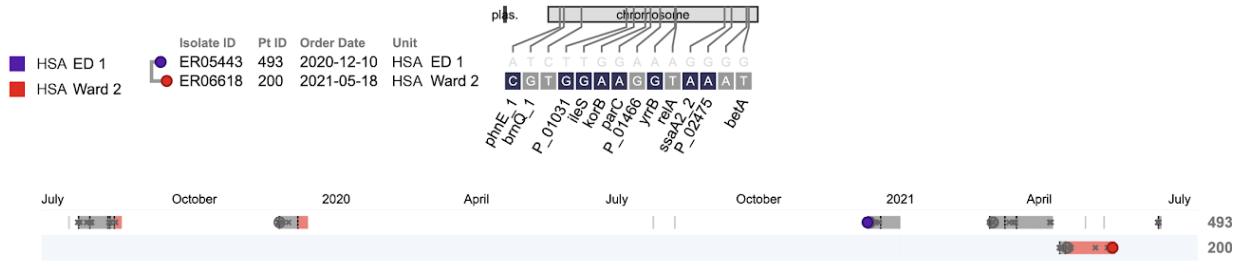
385 References

- 386 Anderson R, Rosenberg A, Garg S, Nahass J, Nenos A, Egorova N, Rowland J, Mari J,
387 LoPachin V. 2019. Establishing the Foundation to Support Health System Quality
388 Improvement: Using a Hand Hygiene Initiative to Define the Process. *J Patient Saf.*
389 <http://dx.doi.org/10.1097/PTS.0000000000000578>.
- 390 Center for Disease Control and Prevention. 2018. Multidrug-Resistant Organism &
391 *Clostridioides difficile* Infection (MDRO/CDI) Module.
- 392 Dupper AC, Sullivan MJ, Chacko KI, Mishkin A, Ciferri B, Kumaresh A, Caban AB, Oussenko I,
393 Beckford C, Zeitouni NE, et al. 2019. Blurred Molecular Epidemiological Lines Between the
394 Two Dominant Methicillin-Resistant *Staphylococcus aureus* Clones. *Open Forum Infectious*
395 *Diseases* **6**. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6735859/> (Accessed September
396 25, 2019).
- 397 Eyre DW, Golubchik T, Gordon NC, Bowden R, Piazza P, Batty EM, Ip CLC, Wilson DJ, Didelot
398 X, O'Connor L, et al. 2012. A pilot study of rapid benchtop sequencing of *Staphylococcus*
399 *aureus* and *Clostridium difficile* for outbreak detection and surveillance. *BMJ Open* **2**.
400 <http://dx.doi.org/10.1136/bmjopen-2012-001124>.
- 401 Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, Sagulenko P, Bedford T,
402 Neher RA. 2018. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* **34**:
403 4121–4123.
- 404 Harris SR, Feil EJ, Holden MTG, Quail MA, Nickerson EK, Chantratita N, Gardete S, Tavares A,
405 Day N, Lindsay JA, et al. 2010. Evolution of MRSA during hospital transmission and
406 intercontinental spread. *Science* **327**: 469–474.
- 407 Köser CU, Holden MTG, Ellington MJ, Cartwright EJP, Brown NM, Ogilvy-Stuart AL, Hsu LY,
408 Chewapreecha C, Croucher NJ, Harris SR, et al. 2012. Rapid whole-genome sequencing
409 for investigation of a neonatal MRSA outbreak. *N Engl J Med* **366**: 2267–2275.
- 410 Lees JA, Harris SR, Tonkin-Hill G, Gladstone RA, Lo SW, Weiser JN, Corander J, Bentley SD,
411 Croucher NJ. 2019. Fast and flexible bacterial genomic epidemiology with PopPUNK.
412 *Genome Res* **29**: 304–316.
- 413 Merrer J, Santoli F, Vecchi CA-D, Tran B, De Jonghe B, Outin H. 2000. “Colonization Pressure”
414 and Risk of Acquisition of Methicillin-Resistant *Staphylococcus aureus* in a Medical
415 Intensive Care Unit. *Infect Control Hosp Epidemiol* **21**: 718–723.
- 416 Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, Phillippy AM. 2016.
417 Mash: fast genome and metagenome distance estimation using MinHash. *Genome Biol* **17**:
418 132.
- 419 Otter JA, Yezli S, French GL. 2011. The role played by contaminated surfaces in the
420 transmission of nosocomial pathogens. *Infect Control Hosp Epidemiol* **32**: 687–699.

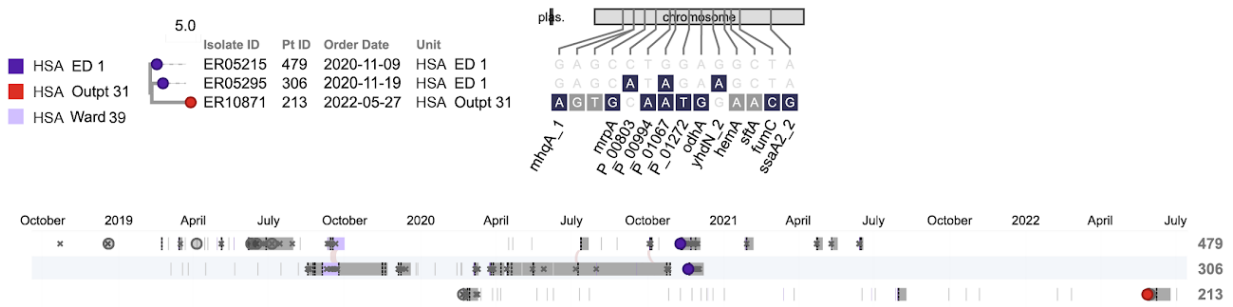
- 421 Otter JA, Yezli S, Salkeld JAG, French GL. 2013. Evidence that contaminated surfaces
422 contribute to the transmission of hospital pathogens and an overview of strategies to
423 address contaminated surfaces in hospital settings. *Am J Infect Control* **41**: S6–11.
- 424 SAS Institute. 2013. *Base SAS 9.4 Procedures Guide: Statistical Procedures, Second Edition*.
- 425 Senn L, Clerc O, Zanetti G, Basset P, Prod'hom G, Gordon NC, Sheppard AE, Crook DW,
426 James R, Thorpe HA, et al. 2016. The Stealthy Superbug: the Role of Asymptomatic
427 Enteric Carriage in Maintaining a Long-Term Hospital Outbreak of ST228
428 Methicillin-Resistant *Staphylococcus aureus*. *MBio* **7**: e02039–15.
- 429 Sullivan MJ, Altman DR, Chacko KI, Ciferri B, Webster E, Pak TR, Deikus G, Lewis-Sandari M,
430 Khan Z, Beckford C, et al. 2019. A Complete Genome Screening Program of Clinical
431 Methicillin-Resistant *Staphylococcus aureus* Isolates Identifies the Origin and Progression
432 of a Neonatal Intensive Care Unit Outbreak. *J Clin Microbiol* **57**.
433 <http://dx.doi.org/10.1128/JCM.01261-19>.
- 434 Treangen TJ, Ondov BD, Koren S, Phillippy AM. 2014. The Harvest suite for rapid core-genome
435 alignment and visualization of thousands of intraspecific microbial genomes. *Genome Biol*
436 **15**: 524.
- 437 Young BC, Golubchik T, Batty EM, Fung R, Lerner-Svensson H, Votintseva AA, Miller RR,
438 Godwin H, Knox K, Everitt RG, et al. 2012. Evolutionary dynamics of *Staphylococcus*
439 *aureus* during progression from carriage to disease. *Proc Natl Acad Sci U S A* **109**:
440 4550–4555.
- 441 Zakhour R, Chaftari A-M, Raad II. 2016. Catheter-related infections in patients with
442 haematological malignancies: novel preventive and therapeutic strategies. *Lancet Infect Dis*
443 **16**: e241–e250.
- 444 Zhou Z, Alikhan N-F, Sergeant MJ, Luhmann N, Vaz C, Francisco AP, Carriço JA, Achtman M.
445 2018. GrapeTree: visualization of core genomic relationships among 100,000 bacterial
446 pathogens. *Genome Res* **28**: 1395–1404.

447 Supplemental Figures

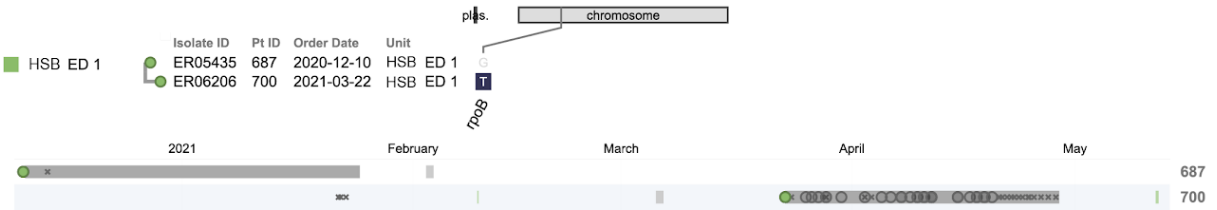
A



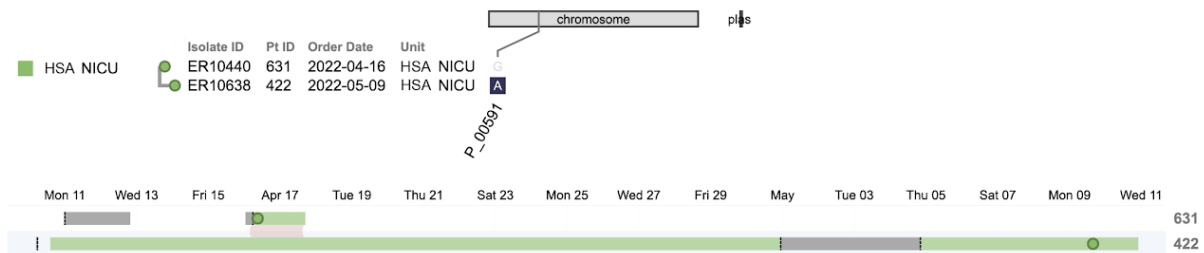
B



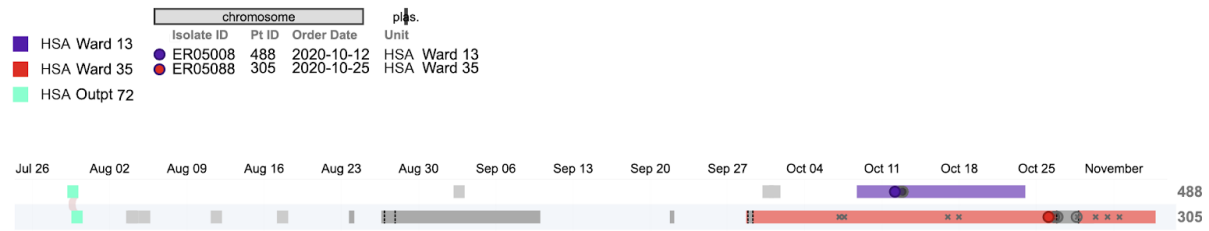
C



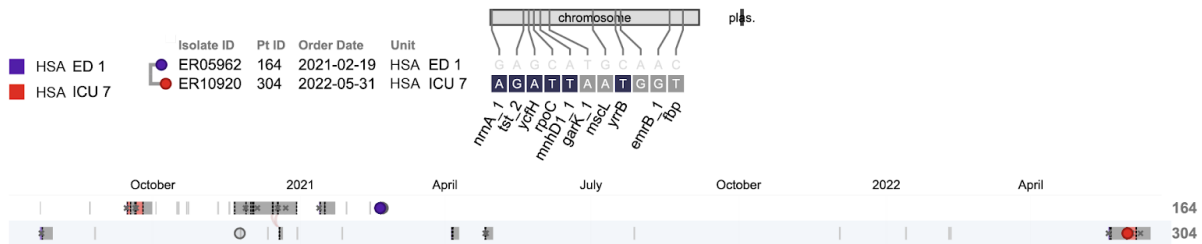
D



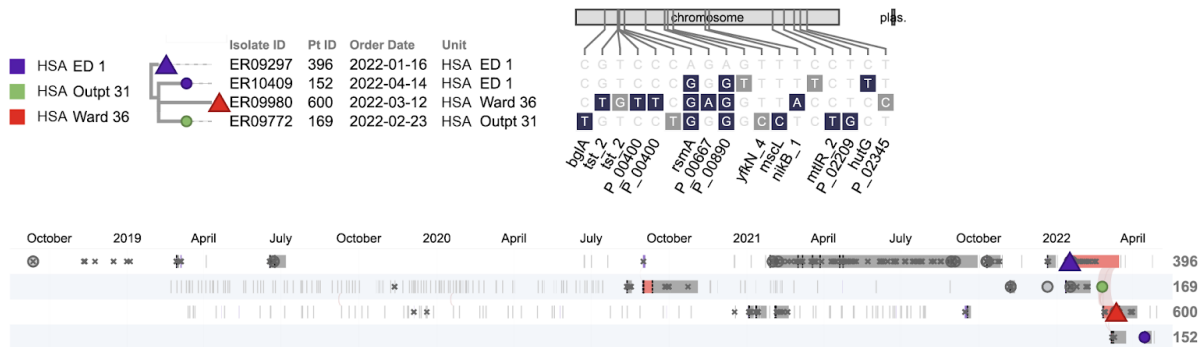
E



F

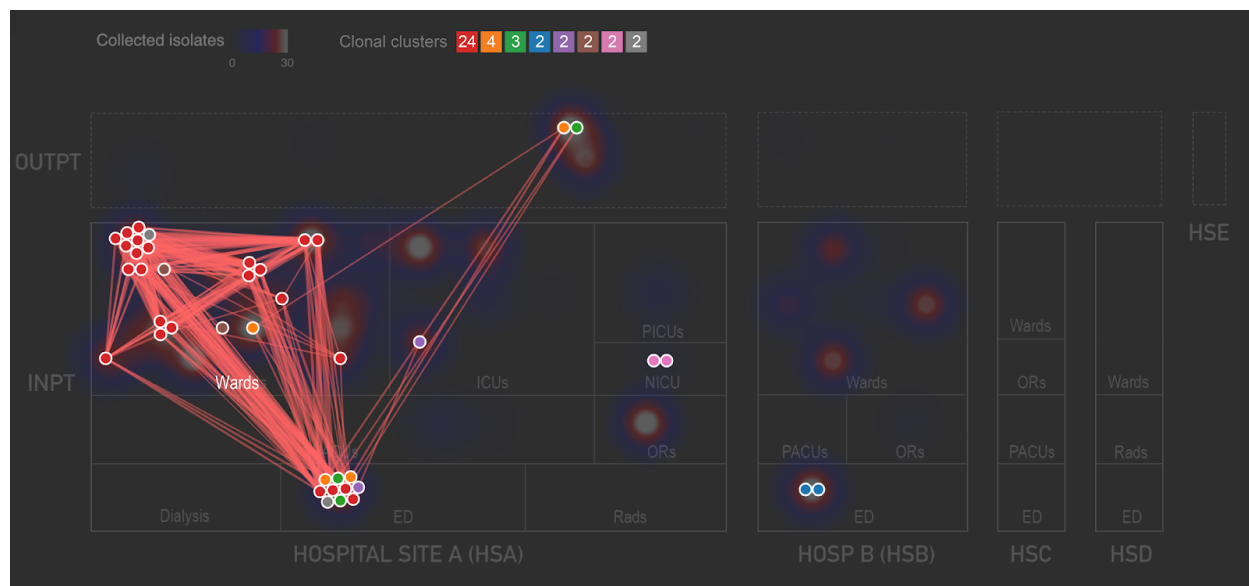


G



448 Figure S1. Other clonal clusters identified by PathoSPOT analysis

449 Phylogenetic trees of core genomes and epidemiologic timelines created by the
 450 “dendro-timeline” visualization in PathoSPOT for the seven smaller clonal clusters identified in
 451 our study (A-G). Layout and drawing conventions are as in Figures 3A and B. All date
 452 information in this figure was recoded to protect health information..



453 Figure S2. Network layout view of clonal clusters

454 The PathoSPOT network layout shows spatial relationships among patients in each of the 8
455 clusters, using the threshold of ≤ 15 SNVs across the entire 24 months of the dataset, and only
456 the first related isolate from each patient. Nodes are colored by clonal cluster, as indicated in the
457 legend at top center, which also lists the numbers of nodes in each cluster. Nodes are laid out
458 spatially by hospital and ward where the patient was sampled, with ward types grouped together
459 (labeled boxes), and nodes from the same ward placed adjacent to one another (the position of
460 each anonymized ward within the box is arbitrary). Genomic links underneath the threshold are
461 depicted as red lines, with many clusters spread across different wards and ward types.
462 Underneath the network map, the total number of positive cultures collected from each location
463 is depicted as a shaded density plot, with the color scale depicted at top left; this highlights
464 heavily sampled locations.