

Complete genome screening of clinical MRSA isolates identifies lineage diversity and provides full resolution of transmission and outbreak events

Mitchell J Sullivan, PhD^{1,2*}, Deena R Altman, MD^{1,3*}, Kieran I Chacko, BS^{1,2}, Brianne Ciferri, MPH^{1,2}, Elizabeth Webster, BS^{1,2}, Theodore R. Pak, BS^{1,2}, Gintaras Deikus, PhD^{1,2}, Martha Lewis-Sandari, BS^{1,2}, Zenab Khan BS^{1,2}, Colleen Beckford, MS^{1,2}, Angela Rendo, BS⁴, Flora Samaroo, BS⁴, Robert Sebra, Ph.D.^{1,2}, Ramona Karam-Howlin, RN³, Tanis Dingle, PhD⁴, Camille Hamula, PhD⁴, Ali Bashir, PhD^{1,2}, Eric Schadt, PhD^{1,2}, Gopi Patel, MD³, Frances Wallach, MD⁵, Andrew Kasarskis, PhD^{1,2}, Kathleen Gibbs, MD^{#6} and Harm van Bakel, PhD^{#1,2*}

1. Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York City, NY, USA
2. Icahn Institute for Data Science and Genomic Technology, Icahn School of Medicine at Mount Sinai, New York City, NY, USA
3. Department of Medicine, Division of Infectious Diseases, Icahn School of Medicine at Mount Sinai, New York City, NY, USA
4. Department of Pathology, Clinical Microbiology, Icahn School of Medicine at Mount Sinai, New York City, NY, USA
5. Department of Medicine, Division of Infectious Diseases, Northwell Long Island Jewish, New York
6. Division of Neonatology and Department of Pediatrics, The Children's Hospital of Philadelphia and The University of Pennsylvania, Philadelphia, PA.

Correspondence to: harm.vanbakel@mssm.edu (HvB)

* These authors contributed equally to this work

These authors are co-senior authors on this work

Abstract

Whole-genome sequencing (WGS) is widely used for studying MRSA evolution but it is not yet a common component of infection prevention practices. We established a continuous genome screening program of all first episode single-patient MRSA blood infections at a major urban hospital. A survey of 132 finished-quality MRSA genomes between 2014-2015 revealed a mixed background of hospital- (USA100/500/800) and community-associated (USA300-NA) lineages, and substantial variation among lineages outside core genomic regions. Megabase-scale inversions caused by homologous recombinations between endogenous prophages were common, particularly among USA100 isolates (28%). We further characterized three transmissions between six adults, and an extended clonal outbreak of USA100 among 3 adults and 18 infants in a neonatal intensive care unit (NICU) lasting 7 months. An analysis of all genetic changes among 23 additional complete isolate genomes obtained during the outbreak provided full spatiotemporal resolution of its origins and progression, which was likely precipitated by ventilator sharing and characterized by several sub-transmissions. The outbreak clone carried distinct mutations in genes with roles in metabolism, drug resistance and persistence. This included a DNA recognition domain recombination in the *hsdS* gene of a Type-I restriction-modification system that altered DNA methylation patterns. RNA-Seq profiling *in vitro* showed that the (epi)genetic changes attenuated *agr* gene expression and upregulated genes involved in stress response and biofilm formation in the outbreak clone, which may have contributed to its spread. Our findings demonstrate the value of a program that provides routine complete genome screening as part of standard infection prevention practices.

Introduction

Hospital-associated (HA) infections with methicillin-resistant *Staphylococcus aureus* (MRSA) are common, impair patient outcomes, and increase healthcare costs (1, 2). MRSA is highly clonal and much of our understanding of its dissemination has relied on lower resolution molecular strain typing methods such as pulsed-field gel electrophoresis (PFGE), *S. aureus* protein A (*spa*) typing, and multilocus sequence typing (MLST) (3), as well as the characterization of accessory genome elements that typify certain lineages and are implicated in their virulence. Examples of the latter include the arginine catabolic mobile element (ACME), *S. aureus* pathogenicity island 5 (SaPI5) and the Panton-Valentine leukocidin (PVL)-carrying ϕ Sa2 prophage in the community-associated (CA) USA300 lineage (4, 5). Molecular typing facilitates rapid screening but has limited resolution to identify transmissions in clonal lineages. Moreover, genetic changes can lead to alteration or loss of typing elements (6–9). As such, WGS has emerged as the gold standard for studying lineage evolution and nosocomial outbreaks (10, 11). Transmission analysis with WGS has been performed largely retrospectively to date (10, 12–15), although it has been used for rapid typing in some cases (16).

In addition to lineage and outbreak analysis, WGS has furthered our understanding of *S. aureus* pathogenicity by providing precise delineation of virulence and drug resistance determinants (17–19), including those related to adaptation to the hospital environment (17, 20). Many of these elements are found in non-conserved ‘accessory’ genome elements that include endogenous prophages, mobile genetic element (MGE), and plasmids (21, 22). The repetitive nature of many of these elements means that they are often fragmented and/or incompletely represented in most WGS studies to date due to limitations of commonly used short-read sequencing technologies, curbing insights into their evolution (22). Recent advances in

throughput of long-read sequencing technologies now enable routine assembly of complete genomes (23, 24) and analysis of core and accessory genome elements (13, 18), including DNA methylation patterns (25), but these technologies have not yet been widely used for prospective MRSA surveillance.

Here we describe the results of a complete genome-based screening program of MRSA blood isolates. During a 16-month period we obtained finished-quality genomes for first blood isolates from all bacteremic patients within one-to-six weeks of culture positivity. In addition to providing detailed contemporary insights into prevailing lineages and genome characteristics, we identified multiple transmission events that had not been recognized based on available epidemiological information. During an outbreak event in the neonatal intensive care unit (NICU) we performed additional sequencing of surveillance and clinical isolates, and were able to provide actionable information that discriminated outbreak-related transmissions and identified multiple sub-transmission events. Further integration of genomic data from outbreak and hospital screening isolates with equipment location tracking data, traced the NICU outbreak origin to adult hospital wards and identified ventilators as a potential vector contributing to the spread of the outbreak between different hospital locations. Finally, comparative genome and gene expression analyses of the outbreak clone to hospital background strains identified genetic and epigenetic changes, including acquisition of accessory genome elements, which may have contributed to the persistence of the outbreak clone.

Results

Complete MRSA genomes reveal genetic diversity in the core and accessory genome

In order to map the genetic diversity of MRSA blood infections at The Mount Sinai Hospital

(MSH) in New York City, US, we sequenced the first positive isolate from all 132 MSH inpatients diagnosed with MRSA bacteremia between fall 2014 and winter 2015. Pacific Biosciences (PacBio) Single molecule real-time (SMRT) long-read length RS-II WGS was used to obtain finished-quality chromosomes for 122 of 132 isolates (92%), along with 145 unique plasmids across isolates (Table S1). We reconstructed a phylogeny from a multi-genome alignment (Fig. 1A, S1A), which identified two major clades corresponding to *S. aureus* clonal complexes 8 (CC8; 45.5% of isolates) and 5 (CC5; 50% of isolates) based on the prevailing multi-locus sequence types (ST) in each clade (ST5 and ST105/ST5, respectively). The CC8 isolates further partitioned among the endemic community-associated (CA) USA300 (80%) and the hospital-associated (HA) USA500 (20%) lineages (Fig. 1B), while CC5 isolates mainly consisted of USA100 (75.8%) and USA800 (15.2%) HA lineages (Fig. 1C). Overall, the phylogeny was consistent with major *S. aureus* lineages found in the NYC region and the US (26). Interestingly, there was considerable variability at loci used for molecular strain typing. Divergence from the dominant *spa* type was apparent in 13.3% of CC8 and 13.6% of CC5 lineage isolates. There were also widespread changes at ACME, PVL, and SaPI5 (Fig. 1B) in USA300 isolates, which are signature elements of this lineage (4, 5); 33.3% either carried inactivating mutations or had partially or completely lost one or more elements (Fig. 1B). Interestingly, we found one PVL-positive USA100 isolate (Fig. 1C) that may have resulted from homologous recombination between a Φ Sa2 and Φ Sa2 PVL prophage (Fig. S1). Thus, complete genomes of MRSA blood isolates demonstrate the mobility of the accessory genome in ways that impact commonly used *S. aureus* lineage definitions.

We further examined larger (>500 bp) structural variants that may be missed by short-read based WGS approaches (13, 14, 27). The multi-genome alignment indicated that between 80.8-88.9% of the sequence in each genome was contained in core syntenic blocks

shared among all 132 genomes (Fig. S2A). Another 9.5-16.8% was contained in accessory blocks found in at least two but not all genomes. Most of these accessory genome elements were lineage-specific and associated with prophage regions and plasmids (Fig. S2B). Finally, 0.8-4.5% of sequence was not found in syntenic blocks and included unique elements gained by individual isolates. For example, a 32 kbp putative integrative conjugative element (ICE) carrying genes encoding proteins involved in heavy metal resistance (cadmium, cobalt, and arsenic) and formaldehyde detoxification was inserted after the *rlmH* gene in the USA800 isolate from p58 (Fig. S2C). Similar arsenic resistance elements have been found in *S. aureus* isolated from poultry litter (28), where its presence has been linked to use of organic arsenic coccidiostats such as roxarsone for growth promotion.

The multi-genome alignment further identified large inversions spanning >250 kbp in 18 genomes (13.6%) (Fig. S2A), which were much more common in CC5 (16 inversions) vs. CC8 lineages (2 inversions) (Fig. 1). The ends of these large inversion events were mainly (94%) located within distinct prophage elements that shared large (>10kb) regions of high sequence similarity (>99%), which meant that the exact cross-over points could not be identified. Notably, 11 inversions spanning ~1.15 Mb occurred between Φ Sa1 and Φ Sa5 in CC5 isolates and could only be resolved by using raw long-read data to phase the small number of variants that uniquely differentiated each prophage (Fig. S3). Other inversions involved cross-overs between prophage pairs of Φ Sa1, Φ Sa3, Φ Sa7, and Φ Sa9. The chimeric prophages that resulted from the inversions consisted of new combinations of the two original prophage elements and contained all genes necessary to produce functional phages based on PHASTER (29) analyses. Taken together, this suggests that prophage elements are key drivers of large inversion events in *S. aureus* that increase prophage diversity.

Identification of transmission events among adults and an outbreak in the NICU

We next compared isolate genomes to identify transmissions between patients. To establish similarity thresholds for complete genomes obtained from long read SMRT sequencing data we first examined baseline single nucleotide variant (SNV) distances between within each lineage. Median pairwise genome differences ranged from 101 SNVs for USA800 to 284 SNVs for USA100 (Fig. S4A). We also examined the extent of divergence among 30 bacteremia isolate pairs collected within a span of one month to 1.4 years from individual patients. Pairwise distances for within-patient isolates were substantially lower than the median for each lineage (Fig. S4B-E), consistent with persistent carriage of the same clone (30, 31), with no more than 10 SNVs separating isolate pairs. Notably, several patients showed variation between isolates collected within a span of several days (Fig. S4B-E); indicative of intra-host genetic diversity. As such, we considered intra-host diversity and genetic drift in aggregate and set a conservative distance of ≤ 7 SNVs to define transmission events in our genome phylogeny. At this threshold we identified one USA300 and three USA100 transmissions involving six adults and three infants (Fig. 1B-C, labeled as T1-T4, Table 1). Complete genome alignments for each event confirmed the absence of structural variants.

In the USA300 transmission case (T1) the presumed index patient p5 was bacteremic with the same clone on two occasions ~3 months apart (3 SNV differences; Fig. S5A). The isolate obtained from the recipient (p33), who was later admitted to the same ward for 7 days at the time of collection, differed from the index strain by only 1 SNV. The USA100 isolates in transmission T2 were collected ~4 months apart and although the patients had overlapping stays, they did not share a ward or other clear epidemiological links (Fig. S5B). In transmission T3, both patients shared a ward for several days (Fig. S5C). The final transmission involving 3

infants (T4) was part of a larger clonal outbreak in the NICU that is discussed in detail below. In summary, although the rate of transmissions in our sample set was low (6.8% of patients), three out of four events had not been recognized based on epidemiological information.

NICU outbreak investigation

Prior to our genomic characterization of the T4 transmission, positive clinical MRSA cultures from three infants in the NICU within five weeks prompted an outbreak investigation and consultation with the New York State Department of Health (NYSDOH). From study day 352 onwards, weekly cultures of composite nasal, umbilicus, and groin surveillance swabs, respiratory cultures for endotracheally intubated patients, and clinical cultures of suspected infected sites were performed for all infants not known to be colonized with MRSA. During four months an additional 41 cultures from 20 infants tested positive for MRSA, bringing the total to 46 isolates from 22 infants. Another 3 positive isolates were obtained from incubators in NICU room 2, from a total of 123 environmental swabs (2.4 %). In collaboration with and following NYSDOH recommendations, surveillance cultures were also obtained from healthcare workers (HCWs) who had provided direct care to newly MRSA-colonized infants, using a concentric circles model. After sampling both nares and both palms, 2 out of 130 (1.5%) HCWs tested were identified as nasally colonized with MRSA strains.

The NYSDOH performed PFGE on 22 isolates, of which 14 patient and three environmental isolates had nearly indistinguishable band patterns (data not shown). This included p90 and p110 in transmission T4 (Fig. 1C; p125 was not tested), indicating that the transmission was part of a larger clonal outbreak. The USA100 (ST105) outbreak clone was resistant to fluoroquinolones, clindamycin, gentamicin and mupirocin, and susceptible to trimethoprim-sulfamethoxazole and doxycycline (Fig. S6, Table S1). This pattern was

uncommon (18.2%) among other USA100 isolates in our study. None of the HCW isolates matched the sequence type or the antibiogram of the outbreak clone (Table S1). Both staff members were successfully decolonized with nasal mupirocin and chlorhexidine gluconate (CHG) baths. While the majority of cases were positive by surveillance, there was morbidity related to the outbreak; five infants developed clinical infections, with three bacteremias, one pneumonia, and one surgical site infection. There were no deaths related to the outbreak.

Complete genomes provide detailed spatiotemporal resolution of NICU outbreak origins and progression

During the outbreak we expanded our genomic screening program to include the first isolate of suspected cases. Four initial isolates were discarded after clinical testing and were not available for genomic analysis. From day 354 onwards we obtained 23 additional complete genomes (Table S1). Of these, 19 genomes from 16 infants and three environmental isolates matched the ST105 outbreak strain type, bringing the total to 22 outbreak genomes from 16 infants and the environment. Four other infants were colonized with MRSA strains that were unrelated to the outbreak and to each other (Table S1). We reconstructed a phylogenetic tree based on core genome alignments of all ST105 isolates in our study, which grouped the 22 NICU isolates with matching antibiograms and/or PFGE patterns in a single clade (Fig. 2A). Surprisingly, this clade also contained 3 MRSA isolates obtained from adult bacteremia patients in other hospital wards prior to the first NICU case. The outbreak clade genomes were ≤ 15 SNVs apart, and the clade as a whole differed from other ST105 isolates by ≥ 41 SNVs. We therefore considered the 3 adult isolates to be part of a larger clonal outbreak that spanned 7 months. Based on the pattern of variants present in ≥ 2 complete genomes we could distinguish 4 distinct outbreak subgroups (Fig. 2A). A core genome-based minimum spanning tree yielded a similar pattern (Fig. 2B).

We next incorporated the available epidemiology and genomic data to reconstruct an outbreak timeline (Fig. 3A). The adult cases were associated with subgroup A that marked the onset of the outbreak had overlapping stays and shared wards, though not always at the same time. Several of the earliest NICU isolates from infants p141, p150, and p151 were not available for genomic analysis (marked X in Fig. 3A), but we note that most cases during this time occurred in rooms 1 and 2. The missing isolate from p141 was susceptible to gentamicin and differed from the PFGE pattern of the outbreak clone by five bands. The other two missing isolates from p150 and 151 matched the antibiogram of the outbreak clone and were therefore considered to be part of the outbreak. The majority of cases were identified during the height of the outbreak on days 357-386 and clustered in subgroup C. All but one of the infants in this subgroup stayed in NICU room 2 before or at the time of culture positivity. The three positive environmental isolates were also obtained from this room, suggesting that a local bioburden lead to a high volume of infants colonized in a short period of time.

The rapid increase in new cases prompted a terminal clean (TC) of the NICU on day 395. As part of this process, all infants were temporarily transferred from the NICU to two different locations in the hospital. During the TC, infant p148 who was colonized with the outbreak clone was placed across the hall from p141 in the pediatric intensive care unit (PICU). A positive surveillance culture in the same subgroup (B) as the p141 was obtained for p148 shortly afterwards (Fig. 3A), suggesting that a transmission had occurred in the PICU. In the following weeks, new positive surveillance cultures were also found for three additional infants (p142, p143 and p146). Each had been admitted after the TC and stayed in room 3 before or at the time of culture positivity. Their isolates comprised subgroup D, suggesting that the outbreak clone spread to this location from the closely related subgroup C linked to room 2 (Fig. 2B, 3A). Thus, each outbreak subgroup (A-D) was associated with a specific location (adult wards, PICU,

and NICU rooms 2 and 3, respectively), indicating that shared locations and environmental bioburden were a dominant factor in the spread of the outbreak clone.

Construction in the NICU and a resulting disruption of infection prevention practices was believed to play a role in the initial transmissions of MRSA. The continued transmission prompted *in situ* simulation and a second TC (Fig. 3A). The simulation efforts reinforced the importance of compliance to infection prevention strategies, patient cohorting, enhanced environmental disinfection, and limiting patient census to decrease bioburden (32). Only one new case (p124) was detected after the second TC. Infant p124 was located the PICU at the time of detection and based on the genomic profile (subgroup C) and earlier positive isolates, the transmission was believed to have occurred prior to the final TC and *in situ* simulation. As such, the workflow improvements were effective in halting the outbreak. The weekly surveillance cultures ended after three consecutive weeks of negative cultures on day 452. The last colonized patient was discharged two months later, and we did not detect the outbreak clone in our hospital-wide genomic screening program in the subsequent two years.

Role of ventilator sharing in the NICU outbreak

While location sharing was a major factor within each outbreak subgroup, it did not explain its spread between different locations, or the link between adult and pediatric cases, which were housed in different buildings and cared for by different HCWs. We focused on a potential role of ventilators in the outbreak based on the observations that: *i*) all NICU outbreak cases were on invasive or non-invasive ventilator support prior to culture positivity; *ii*) the three adult patients were ventilated for at least part of their hospitalizations; and *iii*) prior to identification of the NICU outbreak ventilators were shared between adult and pediatric wards. Ventilator exchange between units was discontinued after the first NICU cases were identified.

The ventilators examined during environmental surveillance tested negative for MRSA, but we could not rule out earlier contamination or contributions of other ventilators. We therefore analyzed equipment usage logs and tracking data provided by the hospital's real-time location system (RTLS), and correlated ventilator movements with specific outbreak subgroups. Six ventilators were shared between NICU outbreak cases, or between adult wards and the NICU (Fig. 3B, numbered 1-6). Within the NICU, the sequential use of ventilator 6 by patient p90 and p133, the timing of their respective culture positivity, and the similarity of their isolate genomes, all supported a role for this ventilator in the transmission to p133. Notably, p133 was the only infant in subgroup C that did not stay in room 2. Likewise, ventilators 2 and/or 5 may have been a factor in the spread from room 2 (subgroup C) to room 3 (subgroup D), especially considering that both rooms were terminally cleaned just prior to the transmission. Ventilator 5 may also have been a transmission vector from p150 to p110. Ventilator 3 was used by p141 and later by p149; however, it is unclear if it played a role in the outbreak, as the first two isolates obtained from p141 after ventilator 3 exposure did not match the outbreak.

Three ventilators (1, 2 and 4) were shared between adult and NICU cases, of which two were moved into NICU prior to the first NICU outbreak case (Fig. 3B). Ventilator 1 was briefly used by adult p64 and then transferred between several locations before it was moved into the NICU and later used by infant p150. The first NICU isolate that matched the outbreak clone by antibiogram was isolated from this patient soon after (Fig. 3A, B). Ventilator 4 was used by adult p91 several weeks before this patient developed bacteremia, except for a 2-day period when it was used by infant p151, shortly before the first NICU outbreak case (Fig. 3B). Infant p151 was cared for in the neighboring PICU at this time and remained there until a positive surveillance isolate was obtained. Finally, ventilator 2 was used by adult p64 during two separate hospital visits, but was only transferred after the outbreak had spread to the NICU. Altogether, the

epidemiological and genomic data suggest that ventilator sharing not only played a role in spread of the outbreak from adult wards to the NICU but was also a factor in subsequent sub-transmissions to patients in different NICU rooms.

Mutations in the outbreak clone alter expression of virulence and persistence factors

Given the extended duration of the outbreak we sought to identify genomic features that could have contributed to its persistence. Comparative genome analysis found 42 non-synonymous or deleterious SNVs and indels in the outbreak clone that were not present in any of the ST105 hospital background strains, affecting 35 genes or their promoter regions (Fig. 4A). The products of these genes were primarily involved in nucleotide, amino acid and energy metabolism, as well as environmental signal processing and drug resistance. Several genes encoding cell wall proteins were also affected, including *gatD*, which is involved in amidation of peptidoglycan (33).

Pan-genome analysis with Roary (34) further revealed 71 genes exclusive to the outbreak strain or infrequently (<33%) present in other MLST105 isolates (Fig. 4B). Most of these genes were associated with three prophage regions and a 43.5 kbp plasmid. The additional genes in prophage A encoded only phage replication or hypothetical proteins. Among the genes in prophage B was an extra copy of *clpB*, which promotes stress tolerance, intracellular replication and biofilm formation (35). Prophage C included an extra copy of the *sep* gene encoding an enterotoxin P-like protein associated with an increased risk of MRSA bacteremia in colonized patients (36). The 43.5 kbp plasmid contained the mupirocin (*mupA*), and gentamicin (*aacA-aphD*) resistance genes (Figure S6B) that explained the distinct susceptibility profile of the outbreak clone. High-level mupirocin resistance (HLR) conferred by *mupA* has been linked to transmissions in previous studies (19, 37, 38). Pan-genome analysis

also revealed a unique variant of the *hsdS* gene in the outbreak strain, which encodes the specificity subunit of a Type I restriction modification (RM) system. Closer examination revealed that a recombination event in one of the DNA recognition sites (Fig. S7) changed a recognition motif from 'CTT' (present at 738 sites, overlapping 595 genes and 120 promoter regions) to 'TGG' (present at 304 sites, overlapping 287 genes and 15 promoter regions), resulting in altered genome-wide ^{6m}A DNA methylation profiles compared to other ST105 isolates (Fig. 4B).

We reasoned that the (epi)genetic changes in the outbreak clone could alter gene expression patterns and provide further insights into the effects of these changes. We therefore compared the gene expression profiles of three representative outbreak isolates (i.e., cases) to the three most similar non-outbreak ST105 strains (i.e., controls) during late-log phase growth. The control strains shared the 43.5 kbp plasmid and most of the prophage elements with the outbreak strain and demonstrated similar growth characteristics (Fig. S8). Differential gene expression analysis showed altered expression of 35 genes (Fig. 4C). Two of these genes were mutated in the outbreak clone; a SNP in promoter region of *sdhC* and a duplication of *clpB*. Methylation changes were found in six genes (17.1%), which was lower than the rate of 27.3% across all genes. Thus, most expression changes appear to be indirect results of (epi)genetic changes. Multiple upregulated genes in the outbreak clone encoded proteins involved in stress and heat shock responses. This included *clpB*, which was increased in copy number in the outbreak vs. control strains, but also *dnaK* and *clpC*, which have been linked to biofilm formation in *S. aureus* and adherence to eukaryotic cells (39, 40). Expression of the gene encoding staphylococcal superantigen-like protein 5 (SSL5) was also increased. SSL5 is known to inhibit leukocyte activation by chemokines and anaphylatoxins (41). Among the downregulated genes, the *agrABC* genes of the accessory gene regulator (*agr*) locus stood out. *Agr* is the major virulence regulator in *S. aureus* (42) and decreased *agr* function in clinical isolates is associated

with attenuated virulence and increased biofilms and surface protein expression (43). Taken together, the nature of the genetic and expression changes in the outbreak clone suggests they may have contributed to its persistence.

Discussion

In this study we implemented a broad WGS screening program at a large quaternary urban medical center, with the aim of tracking circulating clones, to identify transmission events, and to understand the genomic epidemiology of endemic strains impacting human health. To our knowledge, this is the largest set of clinical MRSA isolates from bacteremic patients to undergo long-read sequencing and complete genome assembly to date; significantly increasing the number of finished-quality *S. aureus* genomes. Our results highlight multiple facets of the utility of such data. First, the availability of complete genomes allowed us to precisely map all genetic changes between strains, including prophages, mobile genetic elements, and large genomic inversions. Second, complete reconstruction of outbreak genomes provided additional variation data to map sub-transmission events during a NICU MRSA outbreak. Finally, the combination of genetic and gene expression differences between the NICU outbreak clone and USA100 hospital background revealed genomic features that may have contributed to its persistence.

In our surveillance dataset we identified significant variation in the accessory genome within and between lineages that impacted classical typing schemes for lineage analysis. As such, the stability of HA and CA elements should be considered when using such schemes for lineage analysis. Much of the accessory genome variation occurred in prophage elements, further underscoring their importance in *S. aureus* genome organization (44). We also show that prophages are common drivers of large chromosomal inversions, with evidence of multiple independent events throughout the phylogeny. Inversions were much more frequent in CC5 and

USA100, which may reflect higher similarity between endogenous prophages and/or the increased divergence between isolates in the USA100 lineage. Most inversions could only be resolved by long-read sequencing data and our results, combined with our previous observations among MSSA isolates (45), suggest that prophage-mediated recombinations may be more frequent than previously appreciated. Indeed, one inversion event occurred in the outbreak clone during the spread from the adult wards to the NICU. The impact of genomic inversions on *S. aureus* is unclear and will require further study, but they likely explain the highly chimeric and mosaic structure of *S. aureus* prophages. Notably, non-reciprocal double break-and-join or long gene conversion events can facilitate relevant sequence exchanges between prophages (46). This could lead to the reshuffling of virulence genes and a wider horizontal spread as they become incorporated in phages with different host ranges.

Transmissions were rare during the 16-month surveillance period (6.8% of patients), and, with the exception of the NICU outbreak, each event involved a small number of patients. However, when taking the outbreak into account, transmissions impacted 18.9% of patients we screened, demonstrating that a single large event can have a major impact on infection rates. Complete genome data from our hospital-wide screening program provided key information for outbreak management that could not have been obtained by classical means. First, it provided conclusive differentiation of outbreak from non-outbreak isolates, which helped delineate the final case set and determine when the outbreak ended. While outbreak clone MLST and antibiograms provided a rapid initial screening tool, our surveillance indicated that non-outbreak strains with matching antibiograms were circulating the hospital. Second, analysis of all genetic differences between outbreak cases allowed us to identify sub-transmissions and better understand the chain of events that led to each sub-transmission, resulting in more specificity of intervention. Third, the availability of hospital-wide genomic surveillance data indicated that the

NICU outbreak originated much earlier in unrelated adult wards in a different building and helped identify ventilators as likely transmission vectors.

Complete genome analysis of the outbreak clone revealed a pattern of genetic changes that matched its spatiotemporal spread and aligned with patient locations, suggesting that transmission bottlenecks and local environmental contamination led to a unique genetic signature at each site. Some isolates and isolate subgroups were separated by >10 variants, which is relatively high considering a core genome mutation rate of 2.7-3.3 mutations per Mb per year (14, 30). This suggests that the outbreak may have originated from a genetically heterogeneous source, such as a patient with a history of persistent MRSA colonization that accumulated intra-host variants. It is also possible that the combination of selection pressures and transmission bottlenecks contributed to the diversification of the outbreak clone. Considering all available data, we think the most likely scenario is that the NICU outbreak originated from patient p64 and then spread to other adult patients through direct or indirect contact in shared wards. Ventilator 1, used by adult p64 and infant p150, is the most likely vector for entry into the NICU. Ventilator 4 may have provided a potential second entry route via p151, with subsequent transmissions to p141 and p148 (p151 and p141 had an overlapping stay in the PICU). Such a secondary introduction may explain why the p141 and p148 isolates were more distantly related to all other NICU isolates, but we were not able to test this scenario as the isolates from p151 were no longer available. All subsequent cases could be explained by location relative to other MRSA colonized patients or sharing of MRSA-exposed ventilators.

The outbreak strain differed from the hospital background by multiple mutations of core genes, as well as accessory gene gain and loss. Hundreds of genes were impacted by DNA methylation changes in the gene body or promoter regions, but such genes were depleted rather than enriched among differentially expressed genes. As such, the impact of the

methylation changes on the outbreak clone (if any) was unclear. Nonetheless, a common theme among the genetic and expression changes was the relevance of genes involved in biofilm formation, persistence and quorum sensing. Although the collective impact of the mutations will require further investigation, we speculate that these changes may have contributed an increased persistence of the outbreak clone in the environment.

There are some limitations to our study. Our genomic survey was limited to first positive single-patient bacteremias and transmission rates may be increased when considering non-blood isolates. Moreover, by sequencing single colony isolates we likely did not fully capture intra-host heterogeneity. Although such heterogeneity may be less common among bacteremias, we did encounter variation within some patients which was considered when establishing our transmission thresholds. Finally, while we believe that we have reconstructed the most likely transmission routes and vectors for the NICU outbreak, it is possible that other factors such as spread by HCWs and/or other vectors contributed as well.

In conclusion, we find that the application of complete genome sequencing in the clinical space provides significant benefits for infection prevention and control. In addition to providing contemporary data on the genomic characteristics of circulating lineages, directed intervention and containment of identified transmission events can help prevent further outbreak progression. Although our screening program was limited in scope to bacteremias, early detection of a transmission event between the adult and NICU ward could conceivably have allowed staff to intervene earlier. Completely finished genomes also provide the ability to identify unique elements of particular strain. Accumulating a larger repository of complete and unique genome references and variants associated with successful spreading strains may be key to future outbreak detection and prevention programs by providing high-resolution feature sets for prospective and retrospective data mining purposes.

Materials and methods

Ethics statement

This study was reviewed and approved by the Institutional Review Board of the Icahn School of Medicine at Mount Sinai, and the MSH Pediatric Quality Improvement Committee.

Case review

An investigation of the characteristics of the patients included review of existing medical records for relevant clinical data. Unique ventilator identification numbers and the real-time location system (RTLS) enabled mapping of ventilator locations over time.

Bacterial isolate identification and susceptibility testing.

Isolates were grown and identified as part of standard clinical testing procedures in the Mount Sinai Hospital Clinical Microbiology Laboratory (CML), and stored in tryptic soy broth (TSB) with 15% glycerol at -80°C. VITEK 2 (bioMérieux) automated broth microdilution antibiotic susceptibility profiles were obtained for each isolate according to Clinical and Laboratory Standards Institute (CLSI) 2015 guidelines and reported according to CLSI guidelines (47). Susceptibility to mupirocin was determined by E-test (bioMérieux) and susceptibility to chlorhexidine was tested with discs (Hardy) impregnated with 5 µl of a 20% chlorhexidine gluconate solution (Sigma-Aldrich). Species confirmation was performed with MALDI-TOF (Bruker Biotyper, Bruker Daltonics).

DNA preparation and sequencing

For each isolate, single colonies were selected and grown separately on tryptic soy agar (TSA)

plates with 5% sheep blood (blood agar) (ThermoFisher Scientific) under nonselective conditions. After growth overnight, cells underwent high molecular weight DNA extraction using the Qiagen DNeasy Blood & Tissue Kit (Qiagen, 69504) according to the manufacturer's instructions, with modified lysis conditions. Bacterial cells were lysed by suspending cells in 3 μ L of 100mg/ml RNase A (Ambion, AM2286) and ten μ L of 100 mg/ml lysozyme (Sigma, L1667-1G) for 30 minutes at 37°C, followed by incubation with Proteinase K for one hour at 56°C and two rounds of bead beating of one min each using 0.1mm silica beads (MP Bio) (13). Quality control, DNA quantification, library preparation, and sequencing was performed as described previously (13). Briefly, DNA was gently sheared using Covaris G-tube spin columns into ~20,000 bp fragments, and end-repaired before ligating SMRTbell adapters (Pacific Biosciences). The resulting library was treated with an exonuclease cocktail to remove un-ligated DNA fragments, followed by two additional purification steps with AMPure XP beads (Beckman Coulter) and Blue Pippin (Sage Science) size selection to deplete SMRTbells < 7,000 bp. Libraries were then sequenced using P5 enzyme chemistry on the Pacific Biosciences RS-II platform to >200x genome-wide coverage.

Complete genome assembly and finishing

PacBio SMRT sequencing data were assembled using a custom genome assembly and finishing pipeline (45). Briefly, sequencing data was first assembled with HGAP3 version 2.2.0 (23). Contigs with less than 10x coverage and small contigs that were completely encompassed in larger contigs were removed. Remaining contigs were circularized and reoriented to the origin of replication (*ori*) using Circlator (48), and aligned to the non-redundant nucleotide collection using BLAST+ (49) to identify plasmid sequences. In cases where chromosomes or plasmids did not assemble into complete circularized contigs, manual curation was performed using

Contiguity (50). Genes were annotated using PROKKA (51) and visualized using ChromoZoom (52) and the Integrated Genome Browser (IGB) (53). Interproscan (54) was used to annotate protein domains and GO categories for annotated genes.

Resolution of large genomic inversions

To resolve inversion events catalyzed by two prophage elements (*Staphylococcus phage Sa1* and *Staphylococcus aureus* phage Sa5 with large (>40 kbp) nearly identical regions present in some of the assembled genomes, we developed a phasing approach that took advantage of unique variants present in each element (**Fig. S2A**). Raw (i.e. uncorrected) PacBio reads were first mapped to one of the repeat copies using BWA-MEM (55). Variants were then called with Freebayes (56), and high-quality single nucleotide variants with two distinct alleles of approximately equal read coverage were identified. Analogous to procedures used in haplotype phasing, we then determined which variant alleles were co-located in the same repeat element: if at $\frac{3}{4}$ of the raw reads containing a particular allele also encompassed distinct allele(s) of neighboring variant(s), the alleles were considered linked. In all cases this resulted in two distinct paths through the repeated prophage elements that were each linked to unique sequence flanking each repeat. We then used this information to correct assembly errors and identify *bona fide* inversion events between isolate genomes (**Fig. S2B**). Final verification of corrected assembly was performed by examining the phasing of the raw reads with HaploFlow (57).

Phylogenetic reconstruction and molecular typing

Phylogenetic analyses were based on whole-genome alignments with parsnp (58), using the filter for recombination. The VCF file of all variants identified by parsnp was then used to determine pairwise SNV distances between the core genomes of all strains. For visualization of

the whole-genome alignments, isolate genomes were aligned using sibelia (59) and processed by ChromatiBlocks (<http://github.com/mjsull/chromatiblocks>).

The multi-locus sequence type was determined from whole genome sequences using the RESTful interface to the PubMLST *S. aureus* database (60). Typing of *spa* was performed using a custom script (https://github.com/mjsull/spa_typing). SCCmec typing was done using SCCmecFinder (61). Changes to ACME and SaPI5 were determined using BLASTN and Easyfig. Presence or absence of genes in each locus was determined using BLASTX (62) and a gene was considered to be present if 90% of the reference sequence was aligned with at least 90% identity. Prophage regions were detected using PHASTER. Each region was then aligned to a manually curated database of *S. aureus* phage intergrases using BLASTx to identify their integrase group.

Annotation of antibiotic resistance determinants

Antibiotic resistance gene and variants were annotated by comparing to a manually curated database of 39 known *S. aureus* resistance determinants for 17 antibiotics compiled from literature. BLAST (62) was used to identify the presence of genes in each isolate genome, with sequence identity cutoff $\geq 90\%$ and an e-value cut-off $\leq 1e-10$. Resistance variants were identified by BLAST alignment to the reference sequence of the antibiotic resistance determinant. Only exact matches to variants identified in literature were considered.

Identification of NICU outbreak subgroups

Changes between each outbreak isolate and the p133 reference isolate were identified using GWviz (<https://github.com/mjsull/GWviz>), which uses nucdiff (63) to identify all genomic variants between pairs of strains. Nucdiff in turn uses nucmer to find alignments between two genomes and then identifies large structural rearrangements by looking at the organisation of nucmer

alignments and smaller changes such as SNVs or indels by finding differences between the aligned regions. Briefly, raw PacBio reads were aligned back to each outbreak genome assembly using BWA-MEM (55). Provarvis was then used to detect and associate variants with PROKKA gene annotations, and to determine the number and proportion of raw reads supporting variants in each strain. Variants were selected for further delineation of outbreak subgroups if they were present in two or more isolate genomes and supported by at least ten raw reads in each genome, of which at least 75% confirmed the variant.

A graph of SNV distances between isolates was obtained from a multiple alignment of all outbreak isolates. The minimum spanning tree was then constructed using the minimum spanning tree functionality in the Python library networkX (<https://networkx.github.io/>).

Identification of genetic variants unique to the NICU outbreak clone

To determine SNVs unique to the outbreak isolate the marginal ancestral states of the ST105 isolates were determined using RAxML(64) from a multiple alignment of all ST105s generated using Parsnp. We identified all SNVs that had accumulated from the most recent common ancestor of the outbreak strain and the closest related non-outbreak ST105, and the MRCA of all outbreak strains. SNVs causing nonsynonymous mutations or changes to the promoter region of a gene (defined as <500bp upstream of the start site) were plotted. Orthology was assigned using BLASTkoala (65).

Core and accessory gene content in ST105 outbreak and non-outbreak strains was determined using ROARY. Genes found in more than two outbreak strains and less than 33% of the other ST105 genomes were then plotted along with select methylation data. Phylogenetic reconstruction of ST105 was performed using parsnp and the resulting tree and gene presence

information was visualised using `m.viridis.py` (<https://github.com/mjsull/m.viridis>) which uses the python ETE toolkit (66).

DNA methylation profiling

SMRT raw reads were mapped to the assembled genomes and processed using `smrtanalysis v5.0` (<https://www.pacb.com/products-and-services/analytical-software/smrt-analysis/>).

Interpulse durations (IPDs) were measured and processed as previously described (25, 67) to detect modified N6-methyladenine (^{6m}A) nucleotides.

RNA preparation and sequencing

For RNA extraction, overnight cultures in tryptic soy broth (TSB) were diluted (OD600 of 0.05), grown to late-log (OD600 of ~0.80) in TSB, and stabilized in RNALater (Thermo Fisher). Total RNA was isolated and purified using the RNeasy Mini Kit (Qiagen) according to the manufacturer's instructions, except that two cycles of two-minute bead beating with 1 ml of 0.1 mm silica beads in a mini bead-beater (BioSpec) were used to disrupt cell walls. Isolated RNA was treated with 1 μ L (1 unit) of Baseline Zero DNase (Epicentre) at 37°C for 30 min, followed by ribosomal RNA depletion using the Epicenter Ribo-Zero Magnetic Gold Kit (Illumina), according to the manufacturer's instructions.

RNA quality and quantity was assessed using the Agilent Bioanalyzer and Qubit RNA Broad Range Assay kit (Thermo Fisher), respectively. Barcoded directional RNA-Sequencing libraries were prepared using the TruSeq Stranded Total RNA Sample Preparation kit (Illumina). Libraries were pooled and sequenced on the Illumina HiSeq platform in a 100 bp single-end read run format with six samples per lane.

Differential gene expression analysis

Raw reads were first trimmed by removing Illumina adapter sequences from 3' ends using cutadapt (68) with a minimum match of 32 base pairs and allowing for 15% error rate. Trimmed reads were mapped to the reference genome using Bowtie2 (69), and htseq-count (70) was used to produce strand-specific transcript count summaries. Read counts were then combined into a numeric matrix and used as input for differential gene expression analysis with the Bioconductor EdgeR package (71). Normalization factors were computed on the data matrix using the weighted trimmed mean of M-values (TMM) method (72). Data were fitted to a design matrix containing all sample groups and pairwise comparisons were performed between the groups of interest. P-values were corrected for multiple testing using the Benjamin-Hochberg (BH) method and used to select genes with significant expression differences ($q < 0.05$).

Funding information

This research was supported in part by R01 AI119145 (H.v.B.), the Icahn Institute for Genomics and Data Science (A.K., E.S.), the NIAID-supported NRSA Institutional Research Training Grant for Global Health Research (T32 AI07647), the CTSA/NCATS KL2 Program (KL2TR001435, Icahn School of Medicine at Mount Sinai), and the New York State Department of Health Empire Clinical Research Investigator Program (Aberg, Icahn School of Medicine at Mount Sinai) (D.R.A.) and F30 AI122673 (T.R.P.). The funders had no role in study design, data collection and interpretation, or the decision to submit the work for publication.

Acknowledgements

This research was supported in part through the computational resources and staff expertise provided by the Department of Scientific Computing at the Icahn School of Medicine at Mount Sinai. We thank Drs. Karen Southwick, Eleanor Adams, Lin Ying, and John Kornblum from the New York State Department of Health (NYSDOH) for their consultation during the outbreak and providing PFGE results for a subset of the outbreak isolates.

Author contributions

Methodology: M.J.S, D.A., H.v.B.; Data collection: M.J.S, D.A., K.C., B.C., E.W., T.R.P., G.D., M.L.S., Z.K., C.B., A.R., F.S., K.G., H.v.B.; Data curation: M.J.S, D.A., K.C., B.C.; Analysis: M.J.S, D.A., K.C.; Figures: M.J.S, D.A., K.C., H.v.B.; Writing – original draft: M.J.S, D.A., H.v.B.; Writing – review & editing: All authors; Study design: M.J.S, D.A., H.v.B., K.G.; Supervision: H.v.B., K.G.; Project administration: H.v.B., K.G.; Funding acquisition: H.v.B., D.A., T.R.P., A.K. E.S..

References

1. Klevens RM, Morrison MA, Nadle J, Petit S, Gershman K, Ray S, Harrison LH, Lynfield R, Dumyati G, Townes JM, Craig AS, Zell ER, Fosheim GE, McDougal LK, Carey RB, Fridkin SK, Investigators, Active Bacterial Core surveillance (ABCs) MRSA. 2007. Invasive methicillin-resistant *Staphylococcus aureus* infections in the United States. *JAMA* 298:1763–1771.
2. Grundmann H, Aires-de-Sousa M, Boyce J, Tiemersma E. 2006. Emergence and resurgence of methicillin-resistant *Staphylococcus aureus* as a public-health threat. *Lancet* 368:874–885.
3. DeLeo FR, Chambers HF. 2009. Reemergence of antibiotic-resistant *Staphylococcus aureus* in the genomics era. *J Clin Invest* 119:2464–2474.
4. Diep BA, Gill SR, Chang RF, Phan TH, Chen JH, Davidson MG, Lin F, Lin J, Carleton HA, Mongodin EF, Sensabaugh GF, Perdreau-Remington F. 2006. Complete genome sequence of USA300, an epidemic clone of community-acquired methicillin-resistant *Staphylococcus aureus*. *Lancet* 367:731–739.
5. DeLeo FR, Otto M, Kreiswirth BN, Chambers HF. 2010. Community-associated methicillin-resistant *Staphylococcus aureus*. *Lancet* 375:1557–1568.
6. Glaser P, Martins-Simões P, Villain A, Barbier M, Tristan A, Bouchier C, Ma L, Bes M, Laurent F, Guillemot D, Wirth T, Vandenesch F. 2016. Demography and Intercontinental Spread of the USA300 Community-Acquired Methicillin-Resistant *Staphylococcus aureus* Lineage. *MBio* 7:e02183–15.
7. Montgomery CP, Boyle-Vavra S, Daum RS. 2009. The arginine catabolic mobile element is not associated with enhanced virulence in experimental invasive disease caused by the community-associated methicillin-resistant *Staphylococcus aureus* USA300 genetic

- background. *Infect Immun* 77:2650–2656.
8. Uhlemann A-C, Dordel J, Knox JR, Raven KE, Parkhill J, Holden MTG, Peacock SJ, Lowy FD. 2014. Molecular tracing of the emergence, diversification, and transmission of *S. aureus* sequence type 8 in a New York community. *Proc Natl Acad Sci U S A* 111:6738–6743.
 9. Planet PJ, Diaz L, Kolokotronis S-O, Narechania A, Reyes J, Xing G, Rincon S, Smith H, Panesso D, Ryan C, Smith DP, Guzman M, Zurita J, Sebra R, Deikus G, Nolan RL, Tenover FC, Weinstock GM, Robinson DA, Arias CA. 2015. Parallel Epidemics of Community-Associated Methicillin-Resistant *Staphylococcus aureus* USA300 Infection in North and South America. *J Infect Dis* 212:1874–1882.
 10. Köser CU, Holden MTG, Ellington MJ, Cartwright EJP, Brown NM, Ogilvy-Stuart AL, Hsu LY, Chewapreecha C, Croucher NJ, Harris SR, Sanders M, Enright MC, Dougan G, Bentley SD, Parkhill J, Fraser LJ, Betley JR, Schulz-Trieglaff OB, Smith GP, Peacock SJ. 2012. Rapid whole-genome sequencing for investigation of a neonatal MRSA outbreak. *N Engl J Med* 366:2267–2275.
 11. Price J, Gordon NC, Crook D, Llewelyn M, Paul J. 2013. The usefulness of whole genome sequencing in the management of *Staphylococcus aureus* infections. *Clin Microbiol Infect* 19:784–789.
 12. Azarian T, Cook RL, Johnson JA, Guzman N, McCarter YS, Gomez N, Rathore MH, Morris JG, Salemi M. 2015. Whole-genome sequencing for outbreak investigations of methicillin-resistant *Staphylococcus aureus* in the neonatal intensive care unit: time for routine practice? *Infect Control Hosp Epidemiol* 36:777–785.
 13. Altman DR, Sebra R, Hand J, Attie O, Deikus G, Carpini K, Patel G, Rana M, Arvelakis A, Grewal P, Others. 2014. Transmission of Methicillin-Resistant *Staphylococcus aureus* via

Deceased Donor Liver Transplantation Confirmed by Whole Genome Sequencing. *Am J Transplant* 14:2640–2644.

14. Harris SR, Feil EJ, Holden MTG, Quail MA, Nickerson EK, Chantratita N, Gardete S, Tavares A, Day N, Lindsay JA, Edgeworth JD, de Lencastre H, Parkhill J, Peacock SJ, Bentley SD. 2010. Evolution of MRSA during hospital transmission and intercontinental spread. *Science* 327:469–474.
15. Snitkin ES, Zelazny AM, Thomas PJ, Stock F, NISC Comparative Sequencing Program Group, Henderson DK, Palmore TN, Segre JA. 2012. Tracking a hospital outbreak of carbapenem-resistant *Klebsiella pneumoniae* with whole-genome sequencing. *Sci Transl Med* 4:148ra116.
16. Eyre DW, Golubchik T, Gordon NC, Bowden R, Piazza P, Batty EM, Ip CL, Wilson DJ, Didelot X, O'Connor L, Lay R, Buck D, Kearns AM, Shaw A, Paul J, Wilcox MH, Donnelly PJ, Peto TE, Walker AS, Crook DW. 2012. A pilot study of rapid benchtop sequencing of *Staphylococcus aureus* and *Clostridium difficile* for outbreak detection and surveillance. *BMJ Open* 2.
17. Mwangi MM, Wu SW, Zhou Y, Sieradzki K, de Lencastre H, Richardson P, Bruce D, Rubin E, Myers E, Siggia ED, Tomasz A. 2007. Tracking the in vivo evolution of multidrug resistance in *Staphylococcus aureus* by whole-genome sequencing. *Proc Natl Acad Sci U S A* 104:9451–9456.
18. Benson MA, Ohneck EA, Ryan C, Alonzo F 3rd, Smith H, Narechania A, Kolokotronis S-O, Satola SW, Uhlemann A-C, Sebra R, Deikus G, Shopsin B, Planet PJ, Torres VJ. 2014. Evolution of hypervirulence by a MRSA clone through acquisition of a transposable element. *Mol Microbiol* 93:664–681.
19. Copin R, Sause WE, Fulmer Y, Balasubramanian D, Dyzenhaus S, Ahmed JM, Kumar K,

- Lees J, Stachel A, Fisher JC, Drlica K, Phillips M, Weiser JN, Planet PJ, Uhlemann A-C, Altman DR, Sebra R, van Bakel H, Lighter J, Torres VJ, Shopsin B. 2019. Sequential evolution of virulence and resistance during clonal spread of community-acquired methicillin-resistant *Staphylococcus aureus*. *Proc Natl Acad Sci U S A*.
20. Senn L, Clerc O, Zanetti G, Basset P, Prod'hom G, Gordon NC, Sheppard AE, Crook DW, James R, Thorpe HA, Feil EJ, Blanc DS. 2016. The Stealthy Superbug: the Role of Asymptomatic Enteric Carriage in Maintaining a Long-Term Hospital Outbreak of ST228 Methicillin-Resistant *Staphylococcus aureus*. *MBio* 7:e02039–15.
21. Lindsay JA, Holden MT. 2004. *Staphylococcus aureus*: superbug, super genome? *Trends Microbiol* 12:378–385.
22. Sela U, Euler CW, Correa da Rosa J, Fischetti VA. 2018. Strains of bacterial species induce a greatly varied acute adaptive immune response: The contribution of the accessory genome. *PLoS Pathog* 14:e1006726.
23. Chin CS, Alexander DH, Marks P, Klammer AA, Drake J, Heiner C, Clum A, Copeland A, Huddleston J, Eichler EE, Turner SW, Korlach J. 2013. Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat Methods* 10:563–569.
24. Madoui M-A, Engelen S, Cruaud C, Belser C, Bertrand L, Alberti A, Lemainque A, Wincker P, Aury J-M. 2015. Genome assembly using Nanopore-guided long and error-free DNA reads. *BMC Genomics* 16:327.
25. Fang G, Munera D, Friedman DI, Mandlik A, Chao MC, Banerjee O, Feng Z, Losic B, Mahajan MC, Jabado OJ, Deikus G, Clark TA, Luong K, Murray IA, Davis BM, Keren-Paz A, Chess A, Roberts RJ, Korlach J, Turner SW, Kumar V, Waldor MK, Schadt EE. 2012. Genome-wide mapping of methylated adenine residues in pathogenic *Escherichia coli* using single-molecule real-time sequencing. *Nat Biotechnol* 30:1232–1239.

26. Pardos de la Gandara M, Curry M, Berger J, Burstein D, Della-Latta P, Kopetz V, Quale J, Spitzer E, Tan R, Urban C, Wang G, Whittier S, de Lencastre H, Tomasz A. 2016. MRSA Causing Infections in Hospitals in Greater Metropolitan New York: Major Shift in the Dominant Clonal Type between 1996 and 2014. *PLoS One* 11:e0156924.
27. Copin R, Shopsin B, Torres VJ. 2017. After the deluge: mining *Staphylococcus aureus* genomic data for clinical associations and host-pathogen interactions. *Curr Opin Microbiol* 41:43–50.
28. Williams LE, Detter C, Barry K, Lapidus A, Summers AO. 2006. Facile recovery of individual high-molecular-weight, low-copy-number natural plasmids for genomic sequencing. *Appl Environ Microbiol* 72:4899–4906.
29. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, Wishart DS. 2016. PHASTER: a better, faster version of the PHAST phage search tool. *Nucleic Acids Res* 44:W16–21.
30. Young BC, Golubchik T, Batty EM, Fung R, Lerner-Svensson H, Votintseva AA, Miller RR, Godwin H, Knox K, Everitt RG, Iqbal Z, Rimmer AJ, Cule M, Ip CLC, Didelot X, Harding RM, Donnelly P, Peto TE, Crook DW, Bowden R, Wilson DJ. 2012. Evolutionary dynamics of *Staphylococcus aureus* during progression from carriage to disease. *Proc Natl Acad Sci U S A* 109:4550–4555.
31. Von Eiff C, Becker K, Machka K, Stammer H, Peters G. 2001. Nasal carriage as a source of *Staphylococcus aureus* bacteremia. *N Engl J Med* 344:11–16.
32. Gibbs K, DeMaria S, McKinsey S, Fede A, Harrington A, Hutchison D, Torchen C, Levine A, Goldberg A. 2018. A Novel In Situ Simulation Intervention Used to Mitigate an Outbreak of Methicillin-Resistant *Staphylococcus aureus* in a Neonatal Intensive Care Unit. *J Pediatr* 194:22–27.e5.
33. Figueiredo TA, Sobral RG, Ludovice AM, Almeida JMF de, Bui NK, Vollmer W, de

- Lencastre H, Tomasz A. 2012. Identification of genetic determinants and enzymes involved with the amidation of glutamic acid residues in the peptidoglycan of *Staphylococcus aureus*. *PLoS Pathog* 8:e1002508.
34. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes M, Falush D, Keane JA, Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics* 31:3691–3693.
35. Frees D, Chastanet A, Qazi S, Sørensen K, Hill P, Msadek T, Ingmer H. 2004. Clp ATPases are required for stress tolerance, intracellular replication and biofilm formation in *Staphylococcus aureus*. *Mol Microbiol* 54:1445–1462.
36. Calderwood MS, Desjardins CA, Sakoulas G, Nicol R, Dubois A, Delaney ML, Kleinman K, Cosimi LA, Feldgarden M, Onderdonk AB, Birren BW, Platt R, Huang SS, CDC Prevention Epicenters Program. 2014. Staphylococcal enterotoxin P predicts bacteremia in hospitalized patients colonized with methicillin-resistant *Staphylococcus aureus*. *J Infect Dis* 209:571–577.
37. Hodgson JE, Curnock SP, Dyke KG, Morris R, Sylvester DR, Gross MS. 1994. Molecular characterization of the gene encoding high-level mupirocin resistance in *Staphylococcus aureus* J2870. *Antimicrob Agents Chemother* 38:1205–1208.
38. Udo EE, Jacob LE, Mathew B. 2001. Genetic analysis of methicillin-resistant *Staphylococcus aureus* expressing high- and low-level mupirocin resistance. *J Med Microbiol* 50:909–915.
39. Singh VK, Syring M, Singh A, Singhal K, Dalecki A, Johansson T. 2012. An insight into the significance of the DnaK heat shock system in *Staphylococcus aureus*. *Int J Med Microbiol* 302:242–252.
40. Chatterjee I, Becker P, Grundmeier M, Bischoff M, Somerville GA, Peters G, Sinha B,

- Harraghy N, Proctor RA, Herrmann M. 2005. Staphylococcus aureus ClpC is required for stress resistance, aconitase activity, growth recovery, and death. *J Bacteriol* 187:4488–4496.
41. Bestebroer J, van Kessel KPM, Azouagh H, Walenkamp AM, Boer IGJ, Romijn RA, van Strijp JAG, de Haas CJC. 2009. Staphylococcal SSL5 inhibits leukocyte activation by chemokines and anaphylatoxins. *Blood* 113:328–337.
 42. Novick RP. 2003. Autoinduction and signal transduction in the regulation of staphylococcal virulence. *Mol Microbiol* 48:1429–1449.
 43. Shopsin B, Copin R. 2018. Staphylococcus aureus Adaptation During Infection, p. 431–459. *In* Fong, IW, Shlaes, D, Drlica, K (eds.), *Antimicrobial Resistance in the 21st Century*. Springer International Publishing, Cham.
 44. Xia G, Wolz C. 2014. Phages of Staphylococcus aureus and their impact on host evolution. *Infect Genet Evol* 21:593–601.
 45. Altman DR, Sullivan MJ, Chacko KI, Balasubramanian D, Pak TR, Sause WE, Kumar K, Sebra R, Deikus G, Attie O, Rose H, Lewis M, Fulmer Y, Bashir A, Kasarskis A, Schadt EE, Richardson AR, Torres VJ, Shopsin B, van Bakel H. 2018. Genome plasticity of agr-defective Staphylococcus aureus during clinical infection. *Infect Immun*.
 46. Fortier L-C, Sekulovic O. 2013. Importance of prophages to evolution and virulence of bacterial pathogens. *Virulence* 4:354–365.
 47. Wayne PA. 2015. CLSI. Performance Standards for Antimicrobial Susceptibility Testing; Twenty-Fifth Informational Supplement. CLSI Document M100-S25, Clinical and Laboratory Standards Institute.
 48. Hunt M, Silva ND, Otto TD, Parkhill J, Keane JA, Harris SR. 2015. Circlator: automated circularization of genome assemblies using long sequencing reads. *Genome Biol* 16:294.

49. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009. BLAST+: architecture and applications. *BMC Bioinformatics* 10:421.
50. Sullivan MJ, Ben Zakour NL, Forde BM, Stanton-Cook M, Beatson SA. 2015. Contiguity: Contig adjacency graph construction and visualisation. *PeerJ PrePrints* 3:e1273.
51. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30:2068–2069.
52. Pak TR, Roth FP. 2013. ChromoZoom: a flexible, fluid, web-based genome browser. *Bioinformatics* 29:384–386.
53. Nicol JW, Helt GA, Blanchard SG Jr, Raja A, Loraine AE. 2009. The Integrated Genome Browser: free software for distribution and exploration of genome-scale datasets. *Bioinformatics* 25:2730–2731.
54. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-Y, Lopez R, Hunter S. 2014. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 30:1236–1240.
55. Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv [q-bioGN]*.
56. Garrison E, Marth G. 2012. Haplotype-based variant detection from short-read sequencing. *arXiv [q-bioGN]*.
57. Bachmann NL, Sullivan MJ, Jelocnik M, Myers GSA, Timms P, Polkinghorne A. 2015. Culture-independent genome sequencing of clinical samples reveals an unexpected heterogeneity of infections by *Chlamydia pecorum*. *J Clin Microbiol* 53:1573–1581.
58. Treangen TJ, Ondov BD, Koren S, Phillippy AM. 2014. The Harvest suite for rapid core-genome alignment and visualization of thousands of intraspecific microbial genomes.

Genome Biol 15:524.

59. Minkin I, Patel A, Kolmogorov M, Vyahhi N, Pham S. 2013. Sibelia: A Scalable and Comprehensive Synteny Block Generation Tool for Closely Related Microbial Genomes, p. 215–229. *In Algorithms in Bioinformatics*. Springer Berlin Heidelberg.
60. Jolley KA, Bray JE, Maiden MCJ. 2017. A RESTful application programming interface for the PubMLST molecular typing and genome databases. *Database* 2017.
61. Kaya H, Hasman H, Larsen J, Stegger M, Johannesen TB, Allesøe RL, Lemvigh CK, Aarestrup FM, Lund O, Larsen AR. 2018. SCCmecFinder, a Web-Based Tool for Typing of Staphylococcal Cassette Chromosome mec in *Staphylococcus aureus* Using Whole-Genome Sequence Data. *mSphere* 3.
62. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. *J Mol Biol* 215:403–410.
63. Khelik K, Lagesen K, Sandve GK, Rognes T, Nederbragt AJ. 2017. NucDiff: in-depth characterization and annotation of differences between two sets of DNA sequences. *BMC Bioinformatics* 18:338.
64. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* 30:1312–1313.
65. Kanehisa M, Sato Y, Morishima K. 2016. BlastKOALA and GhostKOALA: KEGG Tools for Functional Characterization of Genome and Metagenome Sequences. *J Mol Biol* 428:726–731.
66. Huerta-Cepas J, Serra F, Bork P. 2016. ETE 3: Reconstruction, Analysis, and Visualization of Phylogenomic Data. *Mol Biol Evol* 33:1635–1638.
67. Flusberg BA, Webster DR, Lee JH, Travers KJ, Olivares EC, Clark TA, Korlach J, Turner SW. 2010. Direct detection of DNA methylation during single-molecule, real-time

- sequencing. *Nat Methods* 7:461–465.
68. Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* 17:10–12.
 69. Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* 9:357–359.
 70. Anders S, Pyl PT, Huber W. 2015. HTSeq--a Python framework to work with high-throughput sequencing data. *Bioinformatics* 31:166–169.
 71. Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics* 26:139–140.
 72. Robinson MD, Oshlack A. 2010. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol* 11:R25.

Tables

Table 1. Clinical characteristics of patients involved in MRSA transmission events

| Event | Patient | Age (years) | Sex | Ward | Study Day | Clinical description |
|-------|---------|-------------|-----|------|-----------|---|
| T1 | p5 | 56 | M | A | 22, 120 | Patient with short gut syndrome complicated by recurrent line infections. Developed a catheter-related MRSA bacteremia. |
| | p33 | 86 | F | A | 119 | Patient with a history of lung adenocarcinoma who was admitted for a fall and vertebral fracture. Developed a peripheral intravenous catheter phlebitis and MRSA bacteremia. |
| T2 | p9 | 57 | M | A | 38 | Patient with end-stage renal disease. Developed a catheter-related MRSA bacteremia. |
| | p48 | 34 | F | Z | 161 | Patient with history of cervical cancer who had percutaneous drain obstruction resulting in MRSA bacteremia. |
| T3 | p118 | 61 | M | B | 435 | Patient with a psychiatric history admitted with polymicrobial bacteremia with MRSA and <i>S. sanguis</i> . |
| | p117 | 31 | M | B | 434 | Patient with a history of sickle cell anemia presented with acute chest syndrome complicated by respiratory failure requiring ventilator support. Developed a catheter-related MRSA bacteremia. |
| T4 | p90 | <1 | F | NICU | 357 | Patient was born at 25 weeks gestation, required ventilator support. Developed bacteremia 11 days after birth. |
| | p110 | <1 | M | NICU | 420 | Patient was born at 25 weeks of gestation, required ventilator support. Had positive surveillance cultures eight days after birth/admission and developed bacteremia 18 days after admission. |
| | p124 | <1 | M | PICU | 455 | Patient was born at 26 weeks of gestation, required ventilator support. |

Figures

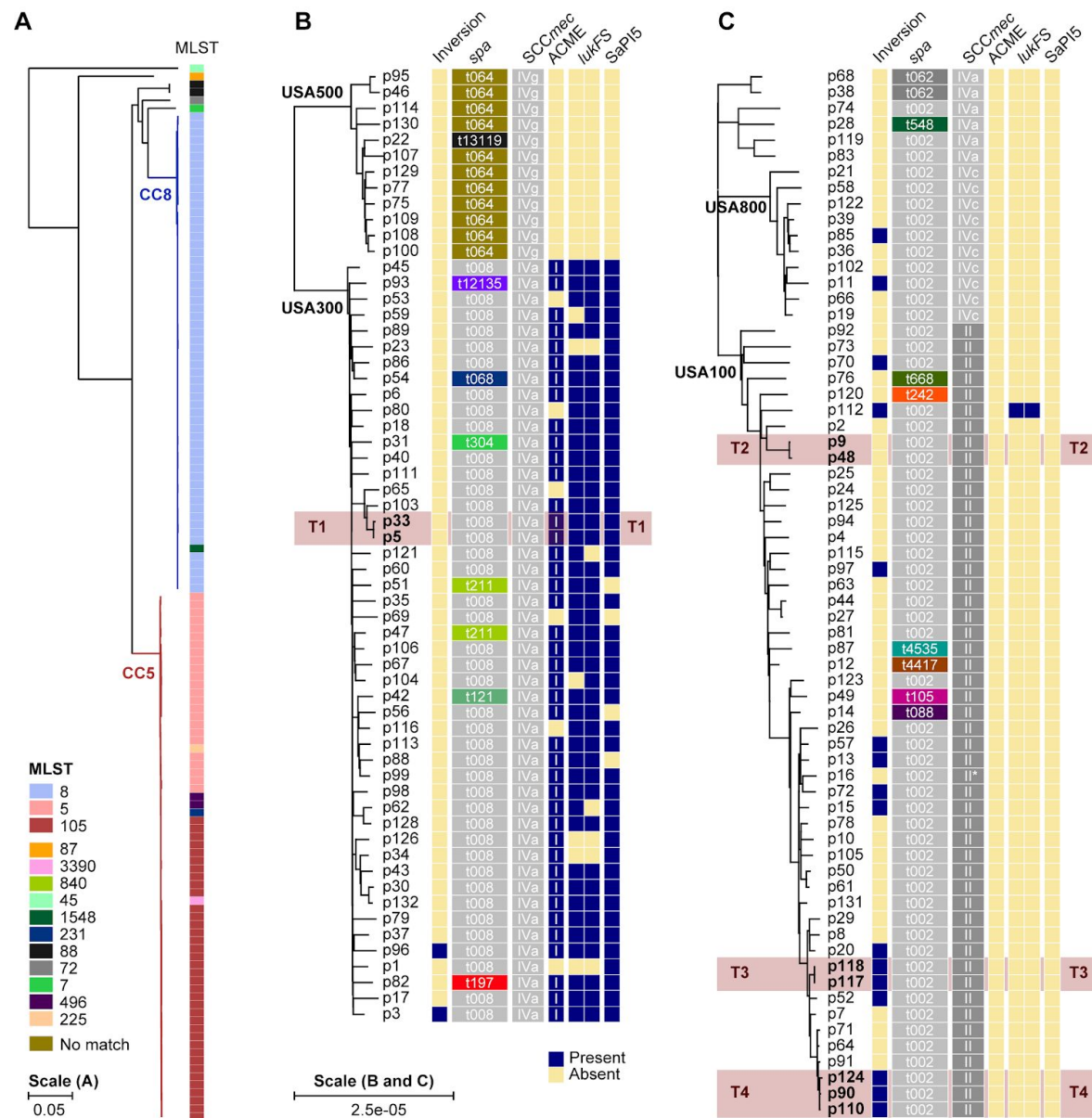


Figure 1. Phylogeny of MRSA bacteremia surveillance isolates.

A) Maximum-likelihood phylogenetic tree based on SNV distances in core genome alignments of 132 primary MRSA bacteremia isolates. CC8 and CC5 clades are shaded in red and blue, respectively. Multilocus sequence types (MLST) for each branch are shown as coloured blocks, with a key at the bottom-left. **B)** Enlarged version of the CC8 clade from A. The isolate identifier is indicated next to each branch, together with blocks denoting the presence of large inversions (>250 kb), *spa* type, *SCCmec* type, and the presence (blue) or absence (yellow) of intact ACME, *lukFS* and SaPI5 loci. The ACME type is indicated in each box. The <

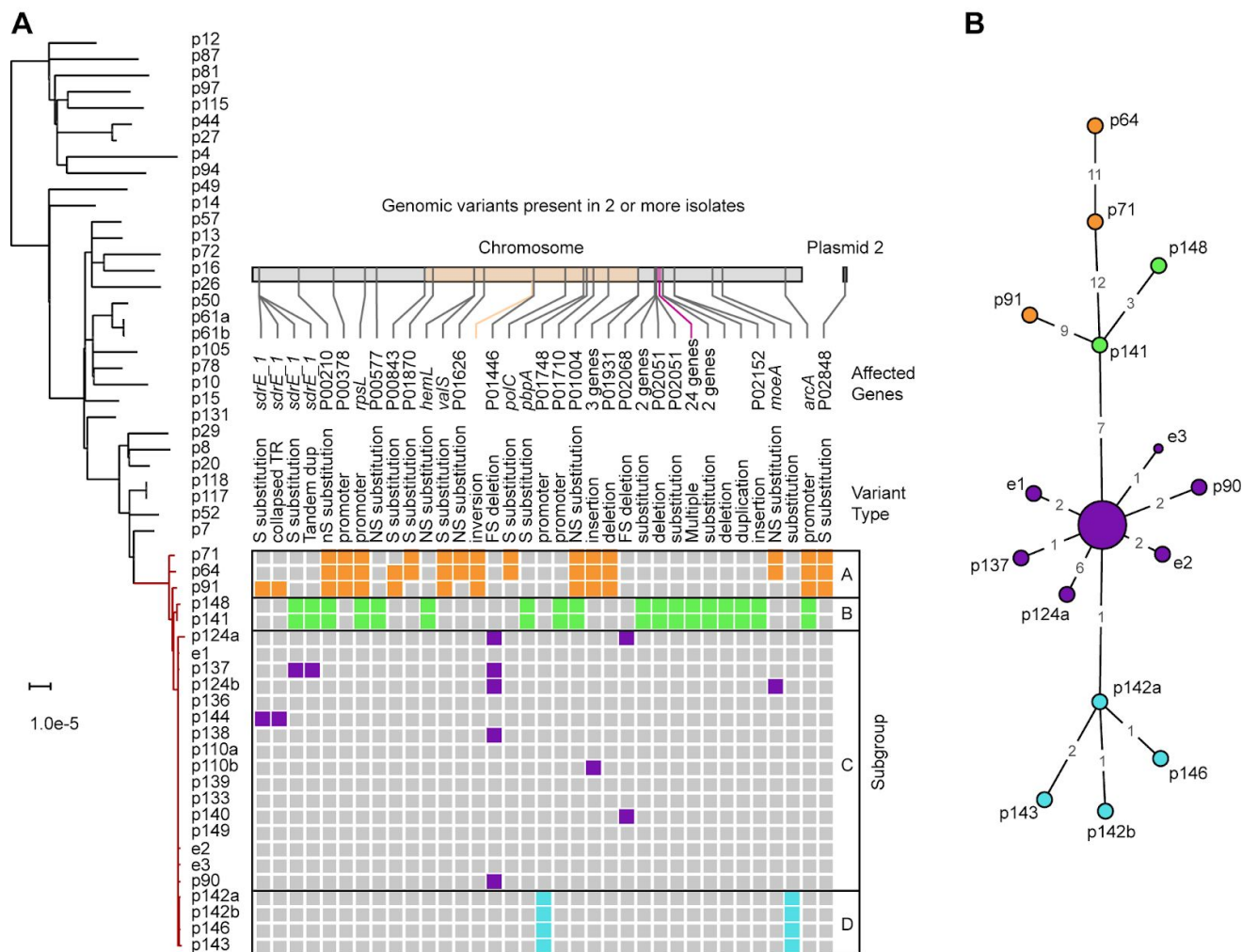


Figure 2. NICU outbreak subgroups and association with adult bacteremia patients.

A) Maximum-likelihood phylogenetic tree based on SNV distances in core genome alignments of 31 ST105 primary bacteremia isolates (black) and 25 outbreak isolates (red). The scale bar indicates the number of substitutions per site. The patient (p) or environmental (e) isolate identifier is shown next to each branch (a/b suffixes indicate multiple isolates from the same patient). Variants present in two or more NICU outbreak isolates, derived from full-length pairwise alignments to the p133 genome, are shown as coloured boxes. Variants are colored according to outbreak subgroups inferred from common variant patterns, as indicated on the right. For each variant the genomic location, affected genes, and type of mutation is shown above the matrix. A 2 Mbp inversion in the adult isolates and a 2,411 bp region containing two substitutions and a deletion in subgroup B is highlighted in the location bar in orange and purple, respectively. **B)** Minimum spanning tree of the 25 outbreak isolates based on SNVs identified in their core genome alignments. The 15 labeled nodes represent individual isolates. The larger central node corresponds to ten isolates with identical core genomes, which includes the p133 reference. Nodes are colored according to the outbreak subgroups shown in panel A. Numbers at edges represent core genome SNV distances.

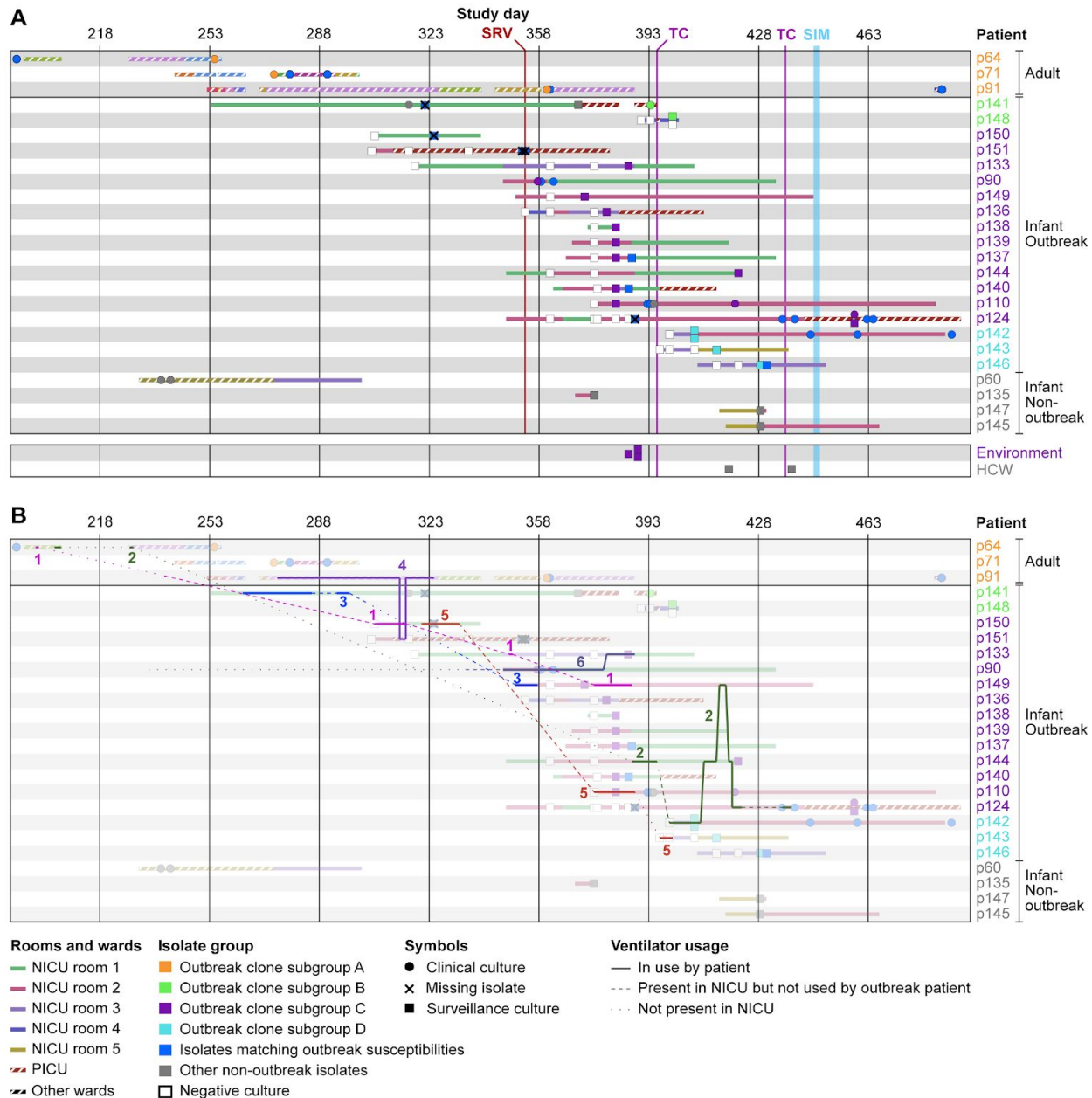


Figure 3. Timeline of the NICU outbreak.

A) Overview of outbreak patient stays and isolates collected during the NICU outbreak. Rows correspond to patients with admission periods shown as horizontal bars. Solid fill patterns denote NICU stays and striped patterns indicate stays in other MSH wards. Fill colors correspond to NICU rooms (solid) or hospital wards (striped). Clinical or surveillance isolates collected during each stay are indicated by symbols, with a key shown below. Patient identifiers and isolate symbols are colored by outbreak subgroup. Timeline scale and key interventions are shown at the top. SRV - start of biweekly surveillance cultures; TC - terminal cleaning; SIM - *in situ* simulation. **B)** Same as A, but with ventilator movements between patients and locations overlaid as lines. Ventilators are numbered and shown in distinct colors. Solid lines correspond to periods that a ventilator was in use by an outbreak patient. Dashed lines indicate when a ventilator was present in the NICU but not used by an outbreak patient. Dotted lines indicate when a ventilator was not in use by an outbreak patient and not present in the NICU. Background colors are muted to facilitate tracking of ventilator movements.

Figure 4. Differentiating features of the NICU outbreak clone compared to the USA100 background.

A) Map of non-synonymous SNVs in genes and promoter regions that are unique to the outbreak clone. Gene identifiers or names are shown next to their genomic location. The SNV type is indicated by colors with a key shown at the top-right. KEGG pathways with two or more genes are indicated on the right (green boxes) and corresponding gene descriptions on the far-right. **B)** Pan-genome analysis of MLST105 isolates showing all genes present in the outbreak clone and absent from at least half of the non-outbreak isolates collected during our study. A maximum-likelihood phylogenetic tree based on SNV distances in core genome alignments is shown on the left with patient (p) or environmental (e) isolate identifiers. Changes in the ^{6m}A methylation profile due to the *hsdS* recombination in the outbreak strain are highlighted in green/blue. Gene presence (yellow) or absence (red) is indicated in a matrix organized by genomic location (top). Gene names and descriptions are shown at the top and bottom of the matrix, respectively. See key on bottom left for more details. **C)** Hierarchical clustering of 35 genes with significant expression differences (FDR $q < 0.05$) between three control and three outbreak strains. Columns correspond to control or outbreak isolates, with labels at the top. Gene names and descriptions are shown on the right. Color shades and intensity represent the difference in normalized log₂ counts per million (CPM) relative to the average gene expression level, with a color key shown below.