# Genomic Diversity of Hospital-Acquired Infections Revealed through Prospective Whole-Genome Sequencing-Based Surveillance

Mustapha M. Mustapha,[a,b] Vatsala R. Srinivasa,[a,b] Marissa P. Griffith,[a,b] Shu-Ting Cho,[a] Daniel R. Evans,[a] Kady Waggle,[a,b] Chinelo Ezeonwuka,[a,b] Daniel J. Snyder,[c] Jane W. Marsh,[a,b] Lee H. Harrison,[a,b] Vaughn S. Cooper,[c,d] Daria Van Tyne[a,d]

[a]Division of Infectious Diseases, University of Pittsburgh School of Medicine, Pittsburgh, Pennsylvania, USA

[b]Microbial Genomic Epidemiology Laboratory, Center for Genomic Epidemiology, University of Pittsburgh, Pittsburgh, Pennsylvania, USA

[c]Department of Microbiology and Molecular Genetics, University of Pittsburgh School of Medicine, Pennsylvania, USA

[d]Center for Evolutionary Biology and Medicine, University of Pittsburgh School of Medicine, Pennsylvania, USA

Mustapha M. Mustapha and Vatsala R. Srinivasa contributed equally. Order of co-first author names was determined by level of seniority in the research group.

**ABSTRACT** Healthcare-associated infections (HAIs) cause mortality, morbidity, and waste of health care resources. HAIs are also an important driver of antimicrobial resistance, which is increasing around the world. Beginning in November 2016, we instituted an initiative to detect outbreaks of HAIs using prospective whole-genome sequencing-based surveillance of bacterial pathogens collected from hospitalized patients. Here, we describe the diversity of bacteria sampled from hospitalized patients at a single center, as revealed through systematic analysis of bacterial isolate genomes. We sequenced the genomes of 3,004 bacterial isolates from hospitalized patients collected over a 25-month period. We identified bacteria belonging to 97 distinct species, which were distributed among 14 groups of related species. Within these groups, isolates could be distinguished from one another by both average nucleotide identity (ANI) and principal-component analysis of accessory genes (PCA-A). Core genome genetic distances and rates of evolution varied among species, which has practical implications for defining shared ancestry during outbreaks and for our broader understanding of the origins of bacterial strains and species. Finally, antimicrobial resistance genes and putative mobile genetic elements were frequently observed, and our systematic analysis revealed patterns of occurrence across the different species sampled from our hospital. Overall, this study shows how understanding the population structure of diverse pathogens circulating in a single health care setting can improve the discriminatory power of genomic epidemiology studies and can help define the processes leading to strain and species differentiation.

**IMPORTANCE** Hospitalized patients are at increased risk of becoming infected with antibiotic-resistant organisms. We used whole-genome sequencing to survey and compare over 3,000 clinical bacterial isolates collected from hospitalized patients at a large medical center over a 2-year period. We identified nearly 100 different bacterial species, which we divided into 14 different groups of related species. When we examined how genetic relatedness differed between species, we found that different species were likely evolving at different rates within our hospital. This is significant because the identification of bacterial outbreaks in the hospital currently relies on genetic similarity cutoffs, which are often applied uniformly across organisms. Finally, we found that antibiotic resistance genes and mobile genetic elements were abundant and were shared among the bacterial isolates we sampled. Overall, this study provides an in-depth view of the genomic diversity and evolutionary processes of

bacteria sampled from hospitalized patients, as well as genetic similarity estimates that can inform hospital outbreak detection and prevention efforts.

**KEYWORDS** whole-genome sequencing, hospital-acquired infections, pangenome, antimicrobial resistance, horizontal gene transfer, bacterial evolution

Healthcare-associated infections (HAIs) affect over half a million people in the United States each year, and annual direct hospital costs for treating HAIs are estimated at over $30 billion (1–3). A relatively small number of bacterial species account for the majority of the burden of antibiotic-resistant HAIs. Organisms belonging to the ESKAPE (*Enterococcus faecium, Staphylococcus aureus, Klebsiella pneumoniae, Acinetobacter baumannii, Pseudomonas aeruginosa,* and *Enterobacter* spp.) group of pathogens are particularly problematic, due to their high burden of HAIs and frequent multidrug resistance (2, 4). In addition, while *Clostridioides difficile* is not highly antibiotic resistant, toxin-producing *C. difficile* lineages associated with significant patient morbidity and mortality have emerged in recent years, making this organism an urgent health threat (5).
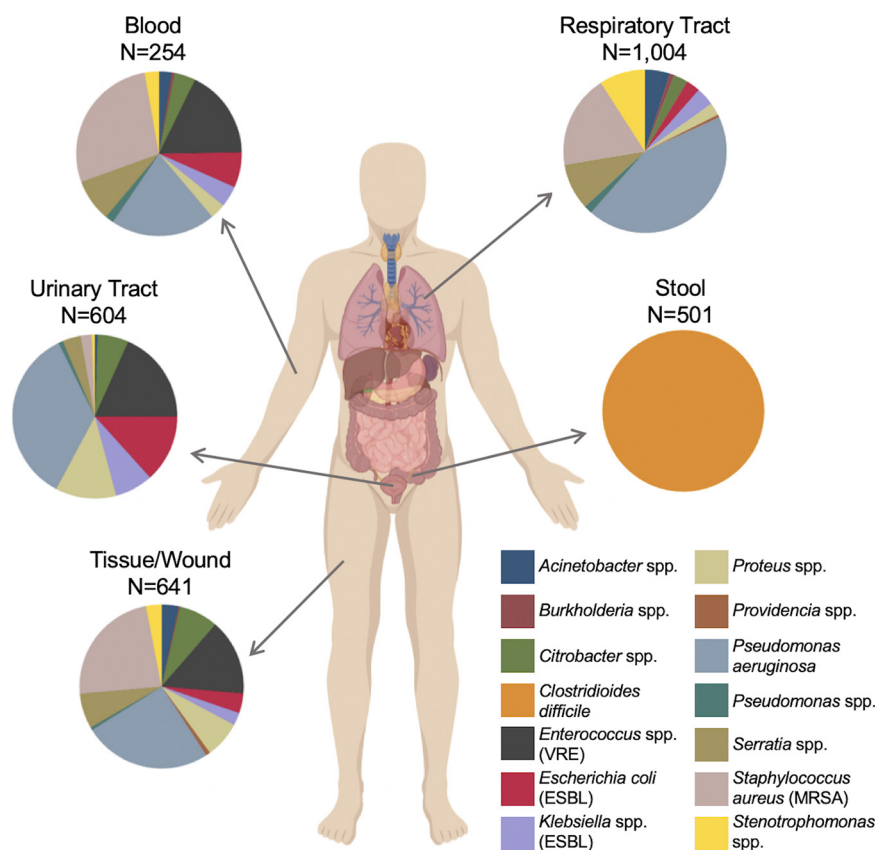
Healthcare institutions such as hospitals and long-term-care facilities constitute a unique ecological niche for the proliferation and spread of antibiotic-resistant pathogens. The hospital environment has a constant flow of vulnerable populations, and widespread use of antimicrobial medications and cleaning agents provide selective pressure for the emergence and expansion of drug-resistant bacteria (6). Likewise, pathogens causing HAIs possess several common biological traits that facilitate their survival and spread in health care environments. These traits include frequent presence and acquisition of antimicrobial resistance, asymptomatic carriage, and the ability to survive for prolonged periods on environmental surfaces such as medical equipment or in water systems (7–9). These factors make health care settings a key contributor to the increase of antibiotic-resistant bacterial infections worldwide.

Epidemiologic surveillance of HAIs requires timely and accurate ascertainment of strain type to identify patients infected with genetically related strains of the same pathogen. Surveillance using whole-genome sequencing (WGS) is the gold standard for the detection of outbreaks and has provided significant insight into the population structure of hospital-associated bacterial infections (10, 11). To improve the detection of hospital-associated transmission at our medical center, we began conducting prospective WGS surveillance of clinical bacterial isolates from hospitalized patients in November 2016, with the aim of identifying previously undetected outbreaks and characterizing pathogen transmission routes. Our approach, called Enhanced Detection of Hospital-Associated Transmission (EDS-HAT), combines prospective bacterial WGS surveillance with data mining of the electronic health record to identify outbreaks and their transmission routes, many of which would otherwise go undetected (12–16). In conducting this work, we have collected and sequenced the genomes of thousands of bacterial isolates. Systematic analysis of these genomes can increase our understanding of the diversity of bacteria causing HAIs (17).

Here, we describe the genomic diversity, evolutionary rates, antimicrobial resistance gene repertoires, and mobile genetic elements carried by over 3,000 bacterial isolates sampled from patients at an academic medical center over 25 months. We uncovered a large and diverse number of species causing HAIs at our center and showed how different population structures and evolutionary rates among these species can impact epidemiologic studies. These findings also shed light on the processes giving rise to bacterial lineages and species. Systematic analyses of antimicrobial resistance genes and mobile genetic elements revealed both species-specific differences as well as broader trends and uncovered new avenues for further investigation.
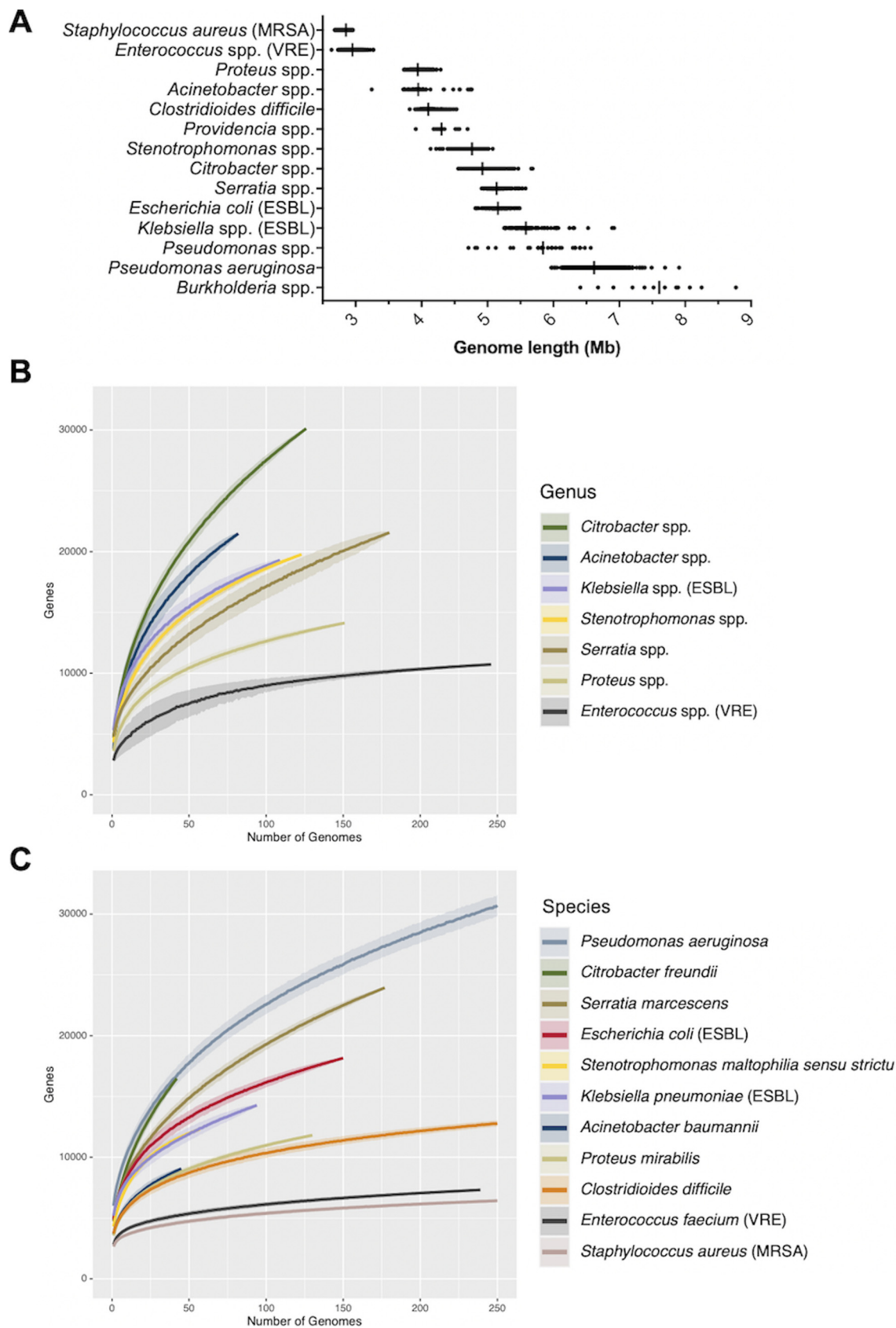
## RESULTS

**Pangenome analysis highlights the diversity of bacteria causing HAIs.** The objective of this study was to use WGS to examine the genetic diversity of HAIs at a single medical center over a multiyear period, and to understand how this diversity

**FIG 1** Species and body site distribution of 3,004 clinical bacterial isolates from hospitalized patients. Isolates were collected from a single hospital over 25 months as part of the Enhanced Detection System for Healthcare-Associated Transmission (EDS-HAT) project. Pie charts show the distribution of isolates belonging to 14 different species groups collected from different types of clinical specimens.

impacts genomic epidemiology and outbreak investigations. A total of 3,004 bacterial isolates collected from 2,046 unique patients at the University of Pittsburgh Medical Center (UPMC) from November 2016 through November 2018 were sequenced and analyzed. Isolates were distributed among 14 groups of species belonging to the same genus, which we called species groups (Tables S1 and S2, Fig. 1). The largest proportion of isolates were sampled from the respiratory tract (33.4%) followed by tissue/ wound (21.3%), urinary tract (20.1%), stool (16.7%, all *C. difficile*), and blood (8.5%) (Fig. 1). The distribution of isolated species was similar between blood and tissue/ wound, while the urinary tract, respiratory tract, and stool samples had different species compositions. *P. aeruginosa* was the most prevalent species isolated, with 863 isolates (28.7% of all isolates) collected from 653 unique patients (Table S1). Other prevalent species included toxin-producing *C. difficile* (16.7%), methicillin-resistant *S. aureus* (MRSA, 14%) and vancomycin-resistant *E. faecalis* and *E. faecium* (VRE, 8.2%). The remaining 10 species groups contained less than 200 isolates each (Table S1). Genome sizes were highly variable, and ranged from a median length of 2.9 Mb for MRSA to 7.6 Mb for *Burkholderia* spp. (Fig. 2A). Pangenome collection curves constructed for genera containing multiple species showed that *Citrobacter* spp. and *Acinetobacter* spp. had the greatest pangenome diversity, perhaps due to the large number of distinct species sampled for these groups (Fig. 2B, Table S2). Pangenome collection curves for individual species showed large differences in pangenome diversity between species (Fig. 2C), with MRSA and VRE *faecium* genomes having the lowest diversity, while *P. aeruginosa*, *C. freundii*, and *S. marcescens* had the greatest pangenome diversity. The large and open pangenome of *P. aeruginosa* is well known (18); however, the pangenome diversity of *C. freundii* and *S. marcescens* are not well described.

FIG 2 Genome length and pangenome size among sampled species. (A) Distribution of genome lengths of isolates belonging to each species group, ordered from shortest to longest median genome length. Vertical lines show median values. (B) Pangenome collection
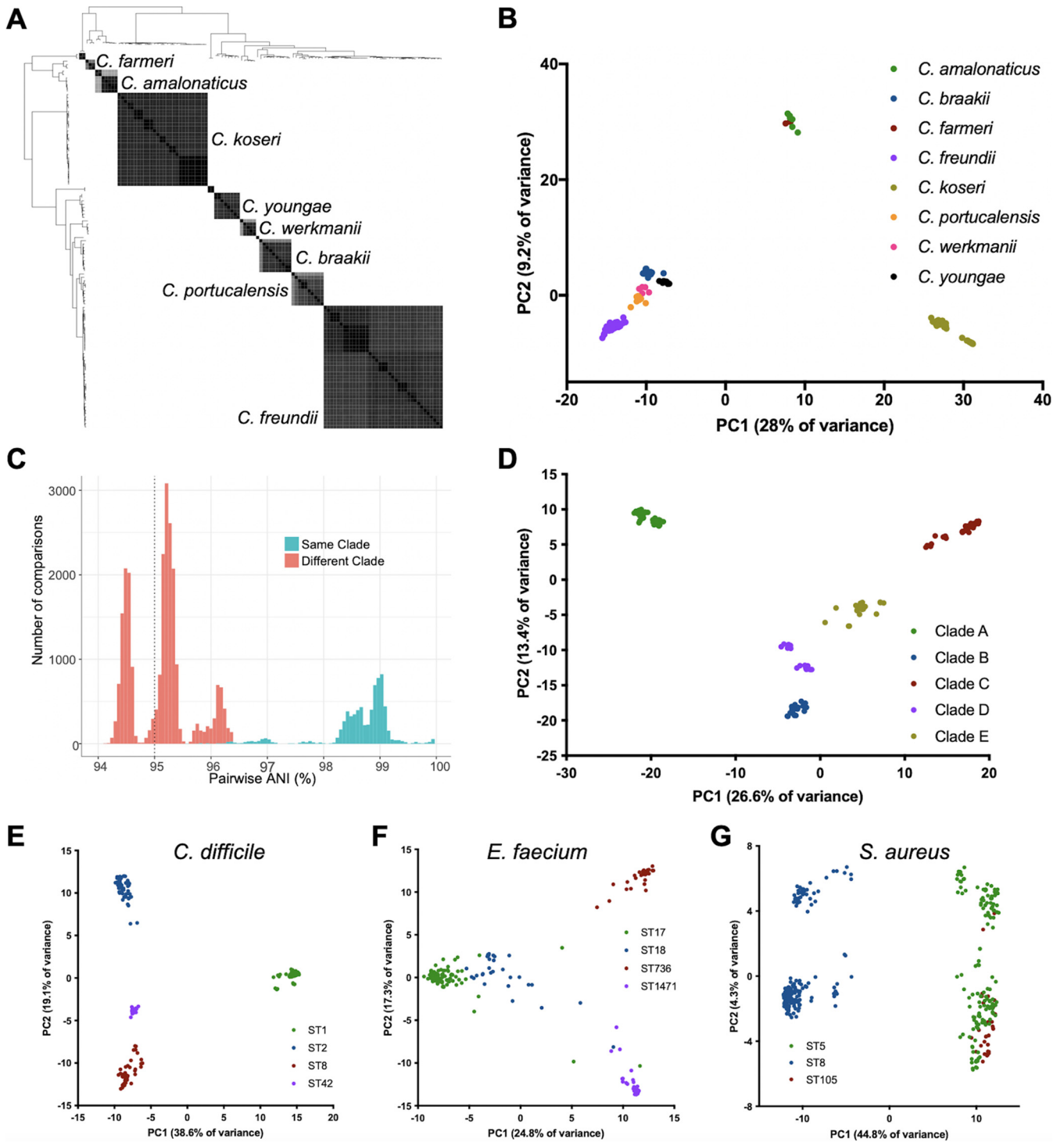
**Differences in bacterial population structures revealed by average nucleotide identity (ANI) and accessory gene content analysis.** Analysis of ANI and accessory genome contents are useful methods for assigning bacterial species, as well as understanding bacterial population structures (19–21). Because the species of each isolate collected by the EDS-HAT project was initially assigned by the clinical microbiology laboratory, we first conducted pairwise comparisons of ANI for all isolate genomes, plus additional reference genomes downloaded from the NCBI database, and used a standard 95% ANI cutoff to group genomes into the same or different species (19). This method resulted in the identification of 97 different species among the collected isolates (Table S2). An example of ANI-based classification of *Citrobacter* spp. is shown in Fig. 3A As expected, several species groups were highly diverse and were composed of multiple different species, including *Acinetobacter* spp., *Burkholderia* spp., *Citrobacter* spp., *Providencia* spp., *Pseudomonas* spp., and *Stenotrophomonas* spp. (Fig. 3A, Fig. S1). Several other species groups, such as ESBL-producing *Klebsiella* spp., *Proteus* spp. and *Serratia* spp., were composed of one dominant species (*K. pneumoniae*, *P. mirabilis*, and *S. marcescens*), and a small number of isolates belonging to other species (Table S1). ANI analysis of *P. aeruginosa* identified 15 isolates (1.7% of all *P. aeruginosa* collected) that could be clearly separated from the rest of the *P. aeruginosa* population by ANI (Fig. S2). These 15 isolates all had greater than 95% ANI with the group 3 PA7 genome (22), indicating that they belonged to this divergent group of *P. aeruginosa*. These results suggest that comparison of ANI-based species identification with information provided by clinical laboratory testing could improve species assignments made for diagnostic purposes.

While ANI measures nucleotide identity in regions that are shared between two genomes, the accessory genes, which by definition are variably present in different genomes, can also be used to differentiate between bacterial species (20, 21). We constructed principal-component analysis plots based on accessory gene content (PCA-A) for species groups containing multiple species (Fig. 3B, Fig. S1). The PCA-A plot for *Citrobacter* spp. isolates was largely congruent with species clustering by ANI (Fig. 3B), and the same was true for *Acinetobacter* spp. and *Stenotrophomonas* spp. as well (Fig. S1). The *S. marcescens* isolates we collected could be clearly separated into five different clades by both ANI and PCA-A; we arbitrarily named these clades A-E (Table S1, Fig. S3). We observed that the pairwise ANI distribution among all *S. marcescens* isolates included comparisons of isolates in different clades that fell below the 95% ANI threshold used to distinguish species from one another (Fig. 3C, Fig. S3). Isolates within each *S. marcescens* clade always shared greater than 95% ANI with isolates in at least one other clade; however, comparisons of isolates in Clade A with isolates in either Clade C or Clade E fell below the 95% ANI threshold for same-species comparisons (Fig. S3). PCA-A clearly separated *S. marcescens* clades from one another (Fig. 3D), suggesting that each clade possessed a unique set of clade-specifying genes (Data set S1). These data suggest that the *S. marcescens* population we sampled may be in the process of diverging into distinct subspecies.

We also explored whether PCA-A could be used to cluster isolates belonging to different genetic lineages within a single species (Fig. 3E to G). We analyzed isolates belonging to the dominant lineages of toxin-producing *C. difficile* (Fig. 3E), VRE *faecium* (Fig. 3F), and MRSA (Fig. 3G), and found in all cases that PCA-A could generally separate isolates belonging to different multilocus sequence types (STs). *C. difficile* isolates belonging to ST1, ST2, ST8, and ST42 were clearly separated from one another (Fig. 3E). *E. faecium* isolates belonging to ST736 were clearly separated from isolates belonging to ST17, ST18, and ST1471, which showed some overlap one another (Fig. 3F). Finally,

**FIG 2** Legend (Continued)
curves for up to 250 genomes from genera containing multiple species and with at least 50 genomes collected. Pangenomes were generated by Roary with an 80% protein identity cutoff. (C) Pangenome collection curves for up to 250 genomes from species with at least 40 genomes collected. Pangenomes were generated by Roary with an 95% protein identity cutoff. Curves show the mean pangenome size and shading shows the standard deviation.

**FIG 3** Average nucleotide identity (ANI) and principal-component analysis of accessory genes (PCA-A) distinguish between and within species. (A) Phylogeny and pairwise ANI values for *Citrobacter* spp. sampled by EDS-HAT. Gray shading indicates ANI values >95%, with darker shading showing higher identity. (B) PCA-A plot for *Citrobacter* species with >2 isolates. (C) Pairwise ANI distribution of *S. marcescens* isolate genomes, showing pairwise ANI comparisons between isolates in different clades that fall below the species cutoff (95% ANI, vertical dashed line). (D) PCA-A plot for *S. marcescens* isolates, showing clear separation of five distinct clades. (E-G) PCA-A plots for dominant sequence types (STs) of *C. difficile* (E), *E. faecium* (F), and *S. aureus* (G).

MRSA isolates belonging to ST8 were clearly separated from isolates belonging to ST5 and ST105; however, the latter STs (which belong to the same clonal complex) were not distinguishable from one another (Fig. 3G). Analysis of gene enrichment among these different STs revealed ST-specific gene repertoires, which were largely composed

**FIG 4** Pairwise SNP distances and genome evolution vary between species. (A) Comparison of within-patient, within-cluster, and between-patient single nucleotide polymorphisms (SNPs) for select species. Pairwise comparisons are shown for all isolate pairs belonging to the same sequence type (ST) within each species. Boxes show the median, 25th and 75th percentiles. (B) Genome evolution rates for dominant STs within *C. difficile* (CD), vancomycin-resistant *E. faecium* (VRE), methicillin-resistant *S. aureus* (MRSA) and *P. aeruginosa* (PSA). Isolates belonging to the four largest STs (three largest for MRSA) of each species were considered, and nucleotide substitution rate (SNPs/genome/year) was calculated for each ST separately. Individual data points are labeled with the corresponding ST, and boxes show the median, 25th and 75th percentiles. (C) Recombination events per mutation (R/Theta) for select species. Each data point represents a distinct ST, and data are grouped by species. STs with at least 10 isolates are shown. Boxes show the median, 25th and 75th percentiles. PRO=*P. mirabilis*, SER=*S. marcescens*, KLP=*K. pneumoniae*, EC=*E. coli*, ACIN=*A. baumannii*.

of predicted mobile element genes and hypothetical proteins (Data set S1). These data suggest that analysis of variable gene content may be a useful complement to SNP-based methods in epidemiologic investigations.

**Genetic diversity and evolutionary rates vary by species.** The EDS-HAT project was designed to detect genetically and epidemiologically connected isolates sampled from different patients, and has successfully identified dozens of clusters containing isolates that share common exposures or transmission chains (14–16, 23). In addition, a significant number of patients in this study were repeatedly sampled. To understand how genetic diversity originated and varied by species, we compared within-patient, within-cluster, and between-patient diversity for six different species. We calculated pairwise core genome SNP distances for all isolate pairs belonging to the same ST, which is commonly used to define clonal lineages within species (Fig. 4A). In all cases,
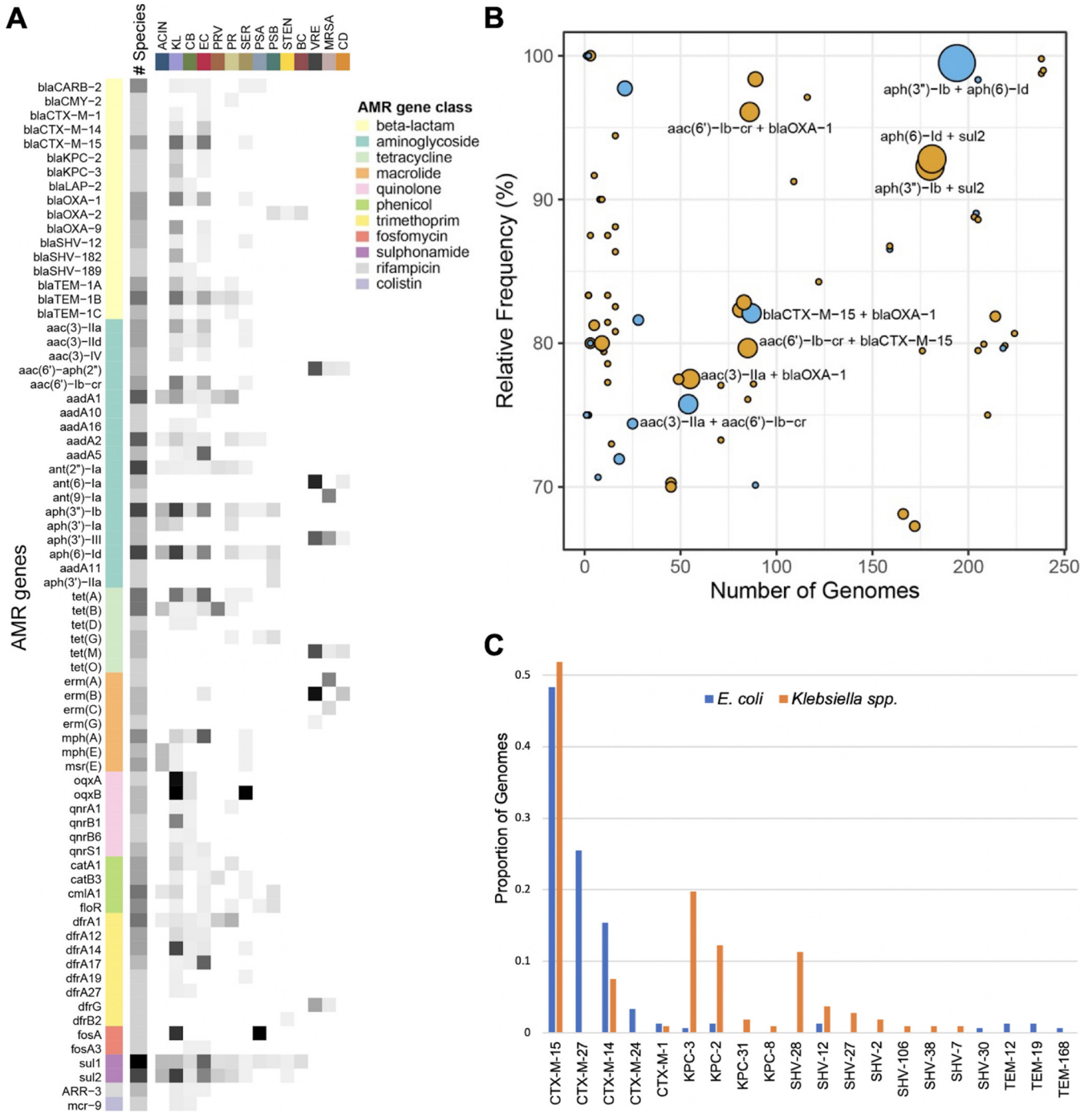
isolates collected from the same patient were on average more similar to one another than those collected from different patients, suggesting that patients were persistently colonized or infected with the same bacterial strain that was repeatedly sampled. Despite only comparing isolates belonging to the same ST, some within-patient comparisons for *P. aeruginosa* resulted in hundreds or thousands of SNPs, which could reflect reinfection with a different strain of a lineage that is diversifying in our hospital, or the presence of hypermutator strains in the infecting bacterial population. Nonetheless, we observed substantial differences in SNP distances among within-patient comparisons across different species, with *C. difficile* isolates having the lowest median pairwise SNPs (2 SNPs), and *P. aeruginosa* having the highest (18 SNPs). These data likely reflect the different genome sizes, evolution and recombination rates, as well as the different biology of these organisms. They also suggest that genomic epidemiologic investigations should consider applying different SNP cutoffs for different organisms.

We recently published a comprehensive analysis of bacterial outbreaks identified in this data set (16). As part of this analysis, we identified clusters of genetically related isolates and then investigated whether the infected patients in each cluster were linked epidemiologically or not. Isolates belonging to epidemiologically linked clusters were highly related to one another, consistent with the rapid spread of clonal outbreak strains (Fig. 4A). Isolates in clusters that did not have epidemiologic links were less closely related to one another (Fig. S4), suggesting possible secondary transmission via patients or other routes that were not sampled.

We next compared the evolutionary rates of the *C. difficile*, VRE, MRSA, and *P. aeruginosa* populations that we sampled. We used TreeTime (24) to estimate the nucleotide substitution rates for the most frequently observed STs for each species (Fig. 4B, Table S3). Consistent with our observations of pairwise SNP differences (Fig. 4A), we found that *C. difficile* had the lowest evolutionary rate, VRE and MRSA had intermediate rates, and *P. aeruginosa* had the highest rate (Fig. 4B). Within each species group, however, we observed a range of nucleotide substitution rates between the different STs that were sampled. Rates overall varied nearly 100-fold among the species and STs we examined, from a minimum of 0.40 SNPs/genome/year for *C. difficile* ST42, to 28.80 SNPs/genome/year for *P. aeruginosa* ST179 (Fig. 4B, Table S3). To understand how recombination might influence these calculations, we used ClonalFrameML (25) to quantify the number of recombination events per point mutation (R/Theta) for each ST across all species for which at least 10 different isolates belonging to the same ST were sampled (Fig. 4C). MRSA genomes were found to have the lowest rates of recombination, while *K. pneumoniae*, *E. coli*, and *A. baumannii* had the highest rates. These data show that rates of nucleotide substitution and recombination are variable across STs as well as across species; this variability should be considered when assessing genomic similarity between isolates during epidemiologic investigations.

**Systematic analysis of acquired antimicrobial resistance (AMR) genes uncovers broad and species-specific trends.** AMR threatens the effective treatment and prevention of bacterial infections. To understand the diversity and distribution of AMR genes among the 3,004 isolates we sampled, we identified acquired resistance genes within each genome by querying the ResFinder database with BLASTn (26) (Fig. S5, Data set S2). The total number of AMR genes identified per genome ranged from 0–19, with an average of 4.6 AMR genes per genome. The species groups carrying the most AMR genes were *Klebsiella* spp. (average 13.1 AMR genes per genome), *E. coli* (7.7 AMR genes per genome), and VRE (average 7.4 AMR genes per genome) (Data set S2). We also classified each AMR gene by drug class and examined the distribution of AMR genes found in more than one species group (Fig. 5A). Several genes encoding aminoglycoside and sulfonamide resistance were observed in the majority of different species groups, suggesting that AMR genes for these antibiotic classes are relatively widespread among bacterial pathogens within our hospital. The Gram-positive species we collected (*C. difficile*, VRE, and MRSA) carried different AMR genes compared to the sampled Gram-negative species, and all Gram-positive species were found to carry the aminoglycoside

FIG 5 Acquired antimicrobial resistance gene abundance and diversity. (A) Prevalence of resistance genes found in more than one species group. Genes are grouped by antibiotic class, and gray shading shows the prevalence of each gene within and across each group. Darker shading indicates higher prevalence. ACIN = *Acinetobacter* spp.; KL = *Klebsiella* spp.; CB = *Citrobacter* spp.; EC = *E. coli*; PRV = *Providencia* spp.; PR = *Proteus* spp.; SER = *Serratia* spp.; PSA = *P. aeruginosa*; PSB = *Pseudomonas* spp.; STEN = *Stenotrophomonas* spp.; BC = *Burkholderia* spp.; VRE = vancomycin-resistant *Enterococcus* spp.; MRSA = methicillin-resistant *S. aureus*; CD = *C. difficile*. (B) Resistance gene co-occurrence. Relative frequency versus number of genomes is plotted for pairs of resistance genes that co-occur at ≥ 50% relative frequency. Blue dots indicate AMR genes in the same drug class, while orange dots indicate genes in different classes. The size of each dot corresponds to the number of different species groups found to carry each pair. AMR gene pairs found in ≥ 4 different species groups are labeled. (C) Distribution of extended-spectrum beta-lactamase (ESBL) and carbapenemase enzymes among *E. coli* and *Klebsiella* spp. isolates.
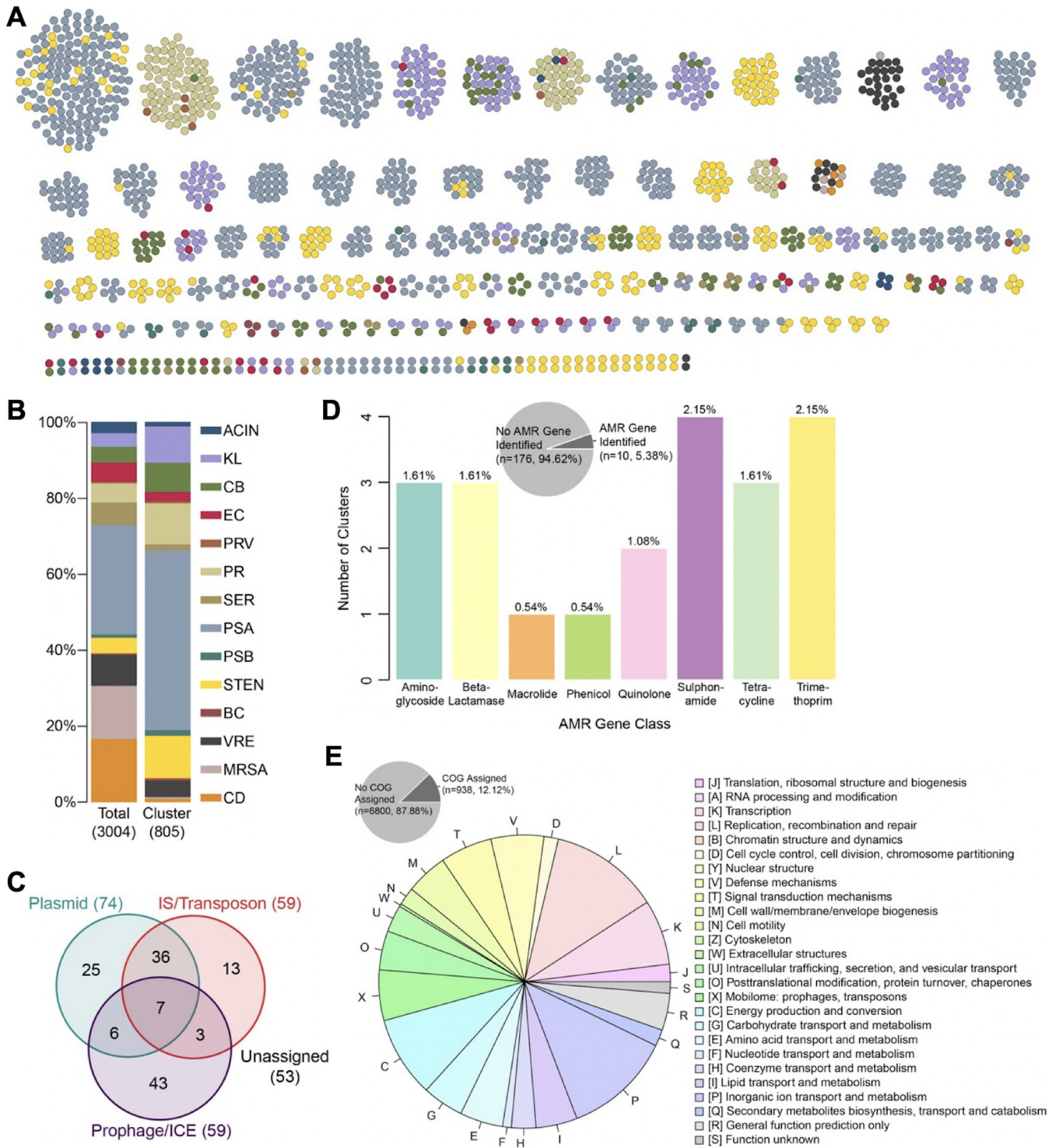
resistance genes *aac(6')-aph(2')* and *aph(3')-III* and the tetracycline resistance gene *tet*(M), albeit at various frequencies (Fig. 5A).

We next examined the co-occurrence of AMR gene pairs across different species groups (Fig. 5B). We found that the aminoglycoside resistance genes *aph(3")-Ib* and

*aph*(6)-*Id* were almost always found together and co-occurred in eight different species groups (all Gram-negative species groups except for *Burkholderia* spp., *Providencia* spp., and *Stenotrophomonas* spp.). Both of these genes also frequently co-occurred with the sulfonamide resistance gene *sul2* (Fig. 5B). A separate aminoglycoside resistance gene, *aac(6')-Ib-cr*, was found to frequently co-occur with the narrow-spectrum beta-lactamase *bla*$_{OXA-1}$, as well as with the extended-spectrum beta-lactamase (ESBL) *bla*$_{CTX-M-15}$. Finally, we examined the distribution of ESBL and carbapenemase enzymes among the ESBL-producing *E. coli* and *Klebsiella* spp. isolates that we sampled (Fig. 5C). The most frequently observed ESBL enzyme was CTX-M-15, which was found in roughly half of all *E. coli* and *Klebsiella* spp. genomes (Fig. 5C). The other half of isolates within each species group carried largely different enzymes from one another, with most *E. coli* isolates carrying other CTX-M-type and a small number of TEM-type ESBLs, while *Klebsiella* spp. isolates carried CTX-M-14 and SHV-type ESBLs. The carbapenemases KPC-2, KPC-3, KPC-8, and KPC-31 were found almost exclusively among *Klebsiella* spp. genomes (Fig. 5C). These data highlight the abundant diversity of AMR genes carried by the bacteria in our hospital. Large surveys of AMR gene occurrence such as this are useful for establishing local patterns of antibiotic resistance, which can guide antibiotic prescribing efforts that are tailored to local pathogen epidemiology.

**Mobile genetic element (MGE) distribution and cargo.** MGEs are frequently found within the genomes of bacteria residing in the hospital environment, and they often encode useful functions like AMR and virulence factors (27). To assess the presence of MGEs in our data set in a systematic and unbiased manner, we used a previously developed approach to identify nucleotide sequences with high homology (>99.9% identity over at least 10Kb) that were present in the genomes of different species (28) (Fig. 6A). This approach resulted in the identification of 186 clusters of shared sequences, which were present in 805 (26.8%) genomes in our data set (Fig. 6B). While each of the 14 different species groups we sampled contained at least one genome encoding a shared sequence, species groups that were particularly enriched for shared sequences included *Klebsiella* spp., *P. aeruginosa*, and *Stenotrophomonas* spp. (Fig. 6B). We next used comparisons with available MGE databases and manual curation to assign an MGE type to each of the 186 clustered sequences based on sequence homology to previously described MGEs (Fig. 6C). We identified similar numbers of sequences that resembled insertion sequences (ISs) or transposons and that resembled prophages or integrative conjugative elements (ICEs). Slightly more sequences showed homology to plasmids, and a large number of sequences resembled multiple MGE types (Fig. 6C). Importantly, 53 (28.5%) shared sequence clusters could not be assigned to an MGE type. Some of these sequences are likely fragments of larger MGEs that lacked genetic features that would enable their classification. Alternately, some of these may constitute novel MGEs.

To understand more about the cargo encoded by the putative MGEs we identified, we first assessed the distribution of AMR genes among the 186 shared sequence clusters we studied (Fig. 6D and Data set S2). Only 10/186 shared sequence clusters (5.4%) carried AMR genes; however, these clusters were found among 116/805 isolates (14.4%). The most frequently observed AMR gene classes (which were each only present in four shared sequence clusters) were sulfonamide and trimethoprim resistance, while aminoglycoside resistance genes, tetracycline resistance genes, and beta-lactamases were each found in three shared sequence clusters (Fig. 6D). We next examined the distribution of clusters of orthologous groups of proteins (COG) categories among all genes present in all shared sequence clusters in our data set. A total of 938 genes (12.1% of all shared sequence cluster genes) had COG categories assigned. Among these genes, the two COG categories observed most frequently were genes involved in replication, recombination and repair, and genes involved in inorganic ion transport and metabolism (Fig. 6E). These data suggest that prominent cargo among the shared sequences we identified included genes for MGE maintenance and mobilization, as well as genes required for the utilization of and resistance to heavy metals, which pathogens frequently encounter in the hospital environment (29).

FIG 6 Mobile genetic element (MGE) distribution and cargo. (A) Clusters of putative MGEs identified in 3,004 study isolate genomes. Nodes within each cluster correspond to bacterial isolates, and are color coded by species group (color key provided in panel B). (B) Distribution of isolates in the entire data set (left) versus isolates encoding one or more putative MGEs (right). (C) Distribution of putative MGEs resembling plasmid, IS/transposon, or prophage/ICE sequences, determined by nucleotide sequence comparisons and manual curation. (D) Distribution of antimicrobial resistance (AMR) genes detected among 186 putative MGEs. (E) Distribution of clusters of orthologous groups of proteins (COG) categories of MGE genes with COG categories assigned.

## DISCUSSION

The broader aim of the EDS-HAT project is to improve the detection of bacterial outbreaks in hospitals, and the project has been successful in this regard (14–16, 23). The EDS-HAT project has also provided a large data set of microbial genomes sampled from thousands of patients within a single medical center over time. Here, we highlight the genetic diversity among bacterial pathogens causing HAIs; understanding this diversity can better inform genomic epidemiology and outbreak investigations. As bacterial WGS becomes increasingly routine in health care settings, this study also provides a baseline for future comparisons, both at our center and elsewhere.

Using comparative genomics methods, we revealed the vast diversity among bacterial pathogens within our hospital. We identified bacteria belonging to 97 different species, which spanned 14 different species groups. We also identified 23 species which have not been previously described, including potentially novel species of *Acinetobacter*, *Citrobacter*, *Proteus*, *Providencia*, *Pseudomonas*, *Serratia,* and *Stenotrophomonas*. A total of 41 isolates (1.4% of sampled isolates) belonged to these potentially novel species, which was a lower proportion than that observed in a prior study of HAIs among ICU patients conducted in 2015 (17). This could be due to additional species having been described in recent years, as well as different inclusion criteria and study populations between the prior study and our own. Further investigation into these new species can aid in the clinical diagnosis of bacteria causing infections.

Our finding that both ANI and PCA-A are effective at distinguishing between different groups at both the species and subspecies levels is consistent with prior studies (30, 31). The 15 *P. aeruginosa* isolates we identified as having 93–94% ANI with the remaining *P. aeruginosa* population is also consistent with prior reports of the *P. aeruginosa* population (32). Conversely, *S. marcescens* is known to have a population structure comprised of multiple clades (33–35), however, we found that pairwise comparisons between some of these clades had less than 95% ANI, suggesting significant divergence and possible ongoing subspeciation. Additionally, we were able to use accessory gene content differences to distinguish between the clades of *S. marcescens* as well as between the dominant genetic lineages of *C. difficile*, VRE *faecium*, and MRSA. Further investigation of these accessory genes would likely enhance our understanding of how different bacterial lineages are able to coexist in the same hospital and could provide useful biomarkers for tracking lineages of interest.

Comparing within-patient versus between-patient genetic diversity can provide guidance in defining SNP cutoffs for outbreak investigations. We found that the number of SNPs among genomes isolated from the same patient at different time points varied by species, with within-patient SNPs being lowest for *C. difficile*, moderate for MRSA and VRE, and greatest for *P. aeruginosa*. Differences between species likely reflect both genome size as well as the biology of these organisms. For example, *C. difficile* can spend long periods of time in a nonreplicative spore state, while *P. aeruginosa* genomes are more than double the size of MRSA and VRE genomes. The SNP distances among same-patient isolates we observed are comparable to those used in outbreak investigations in our setting and elsewhere (14, 36, 37). The EDS-HAT project is currently using SNP cutoffs of ≤2 for *C. difficile* and ≤15 for all other organisms (16). We suggest that researchers performing genomic epidemiology studies should calibrate SNP cutoffs using comparisons of same-patient and confirmed outbreak isolates from their own setting and with their own bioinformatic analysis pipelines.

Evolutionary rates assessed for the four most common species in our hospital were also consistent with previous studies (38, 39). Within each species, however, we observed differences in substitution rates that were 10-fold or more, suggesting that different genetic lineages of the same species may be evolving at different rates within our hospital. While it is possible that these estimates could change with a different or larger sample size, exploring how evolution varies by genetic lineage would provide important insights into the emergence and persistence of these pathogens in the hospital environment. When we compared evolutionary rates across different species, we

observed a nearly 100-fold difference in substitution rates, which was most pronounced when comparing *C. difficile* with *P. aeruginosa*. These differences likely drive the differences in pairwise SNP distances we observed between these species, and further suggest that different SNP cutoffs should be considered for different bacterial species for the purposes of hospital outbreak investigations. In the future we also plan to examine more closely rates of bacterial evolution in the same patients over time.

This study established the diversity of acquired antimicrobial resistance genes among pathogenic bacteria circulating at our hospital and provides a point of comparison with other studies of antibiotic resistance spread in the hospital environment (23, 28, 40, 41). Using a moderately stringent cutoff for sequence coverage and identity, we found that aminoglycoside and sulfonamide resistance genes were highly abundant and were found in the majority of species that we sampled. Although the presence of aminoglycoside resistance is well documented among both Gram-positive and Gram-negative bacteria—and more specifically among the ESKAPE pathogens—less attention has been focused on sulfonamide resistance (42–44). The co-occurrence of *aph (3″)-Ib*, *aph*(6)-*Id*, and *sul2* that we observed has also been previously observed in a variety of different genetic contexts, including in plasmids, integrative conjugative elements, and chromosomal genomic islands (43, 45). Additionally, we found that the ESBL enzyme $bla_{CTX-M-15}$ was widely distributed among both *E. coli* and *Klebsiella* spp. isolates, which is consistent with prior reports (46). Among the other ESBL-producing *E. coli* and *Klebsiella* spp. isolates collected, ESBL enzymes were largely restricted to one species group or the other. Finally, while we did not explicitly collect carbapenemase-producing organisms during this study period, a subset of the ESBL-producing *E. coli* and *Klebsiella* spp. isolates collected also carried carbapenemase enzymes. Co-occurrence of ESBL enzymes and carbapenemases was more frequent among *Klebsiella* spp., especially ST258 *K. pneumoniae* (23).

This study also offers an overview of highly similar sequences (which we suspect largely belong to MGEs) shared among the genomes of distantly related bacteria sampled from patients residing in the same hospital environment. We found that *Enterobacteriaceae* such as *Klebsiella* spp. and *Citrobacter* spp., as well as *P. aeruginosa* and *Stenotrophomonas* spp., were overrepresented among shared sequence clusters compared to their overall distribution in the data set. Most of the shared sequences identified in *Enterobacteriaceae* genomes resembled sequences carried on plasmids, consistent with the frequent plasmid exchange known to happen among species in this family (47). While we did not fully resolve these plasmids in this study, investigating plasmid transmission dynamics in this setting will be a focus of our future work. In contrast to the sequences shared between *Enterobacteriaceae* genomes, many of the shared sequences identified among *P. aeruginosa* and *Stenotrophomonas* spp. resembled prophages and integrated conjugative elements, suggesting that these organisms may rely on different MGEs to exchange genetic material. Somewhat surprisingly, our analysis identified fewer shared sequences carrying AMR genes compared to a prior study we conducted within the same hospital (28). This may be due to our use of a longer sequence length cutoff for shared sequence identification in this study, as AMR genes are known to be carried on smaller MGE units that can rapidly shuffle, interchange, and mutate (48). Finally, we found it notable that genes encoding metal transport and resistance were frequently observed within the shared sequences we identified. Inorganic ions are required for catalysis of many bacterial enzymes (49), and heavy metals such as silver, copper, and mercury have long been used as disinfectants in hospitals (50). Further study of MGEs encoding metal-interacting genes will be a focus of our future work.

This study had several limitations. The organisms we collected were prespecified, and certain groups, such as *Enterobacter* spp. or carbapenemase-producing organisms without a noted ESBL phenotype, were not collected. Furthermore, our definition of "hospital-acquired infections" was quite broad; some of the collected isolates likely represent commensal organisms or pathogen colonization, rather than true infection.

We also cannot say for sure whether the sampled bacteria were acquired from the health care setting or not, as we only considered bacterial isolates from clinical specimens and did not include environmental sampling. Additionally, our 25-month collection window was quite short, thus we were unable to draw conclusions regarding trends over time. We focused our study of AMR genes on acquired resistance genes, and the coverage and identity cutoffs we used may have caused us to miss divergent and potentially novel resistance genes. Finally, the inclusion of both broad species groups as well as more defined sets of specific pathogens made it difficult to conduct systematic analyses or draw broader conclusions across the entire data set. Nonetheless, the large number of isolates collected offers a high-resolution view of the genomic diversity and evolution of important bacterial pathogens found within our hospital. Our future work will include following these bacterial populations over time and comparing our results with similar studies conducted in other settings.

In assessing the genomes of major infection-associated bacterial species isolated from patients at our hospital, we have provided a survey of the genomic diversity of bacterial HAIs at a single clinical center. Our findings demonstrate that studying population dynamics and evolution of these pathogens can inform genomics-based outbreak investigations. In addition to forming a basis for future comparisons, this study also provides a deeper understanding of the breadth of different species that cause HAIs and demonstrates the utility of systematic genome sequencing and comparative genomics analysis of clinical bacterial isolates from hospitalized patients.

## MATERIALS AND METHODS

**Isolate collection.** Bacterial isolates were collected from the University of Pittsburgh Medical Center (UPMC) Presbyterian Hospital, an adult tertiary care hospital with over 750 beds, 150 critical care unit beds, more than 32,000 yearly inpatient admissions, and over 400 solid organ transplants per year. All isolates were collected from hospitalized patients and were isolated from clinical cultures prompted by clinician suspected infection. Clinical specimens were processed by the UPMC clinical microbiology laboratory. Processing varied by sample type but always resulted in isolation of a representative bacterial clone (or clones) from each specimen for the purposes of species identification and antimicrobial susceptibility testing. Isolates were collected from November 2016 through November 2018 from admitted patients as part of a prospective genomic epidemiology surveillance project called Enhanced Detection System for Healthcare-Associated Transmission (EDS-HAT) (16). Inclusion criteria were hospital admission greater than 2 days before the culture date, and/or a recent inpatient or outpatient UPMC hospital encounter in the 30 days before the culture date, and isolation of an organism on a defined list of high-priority health care-associated pathogens. After patients to be included in the study were identified, the bacterial clone that was isolated and tested by the clinical microbiology laboratory was streaked onto a blood agar plate (BD, Franklin Lakes, NJ), grown overnight at 37°C, and then the plate was scraped and genomic DNA was extracted from the bacterial pellet. For *Clostridioides difficile*, stool specimens that were culture-independent diagnostic test–positive for *C. difficile* were cultured to isolate a single representative clone. A total of 3,004 isolates were included in this study (Table S1).

Isolates were classified as belonging to one of 14 groups of related species, which we called species groups: *Acinetobacter* spp., *Burkholderia* spp., *Citrobacter* spp., *Clostridioides difficile*, *Enterococcus* spp. (VRE), *Escherichia coli* (ESBL-producing), *Klebsiella* spp. (ESBL-producing), *Proteus* spp., *Providencia* spp., *Pseudomonas aeruginosa*, *Pseudomonas* spp. (non-*aeruginosa*), *Serratia* spp., *Staphylococcus aureus* (MRSA), and *Stenotrophomonas* spp. Isolate collection was limited to only toxin-producing strains of *Clostridioides difficile*, vancomycin-resistant *Enterococcus* spp. (VRE), extended-spectrum beta-lactamase (ESBL)-producing *Escherichia coli* and *Klebsiella* spp., and methicillin-resistant *Staphylococcus aureus* (MRSA). This study was approved by the University of Pittsburgh Institutional Review Board and was classified as being exempt from patient-informed consent.

**Whole-genome sequencing and genome assembly.** Genomic DNA was extracted from pure overnight cultures of single bacterial colonies using a Qiagen DNeasy Tissue kit according to the manufacturer's instructions (Qiagen, Germantown, MD). Illumina library construction and sequencing were conducted using an Illumina Nextera DNA Sample Prep kit with 150 bp paired-end reads, and libraries were sequenced on the NextSeq 550 sequencing platform (Illumina, San Diego, CA). Selected isolates were re-sequenced with long-read technology on a MinION device (Oxford Nanopore Technologies, Oxford, United Kingdom). Long-read sequencing libraries were prepared and multiplexed using a rapid multiplex barcoding kit (catalog SQK-RBK004) and were sequenced on R9.4.1 flow cells. Base-calling on raw reads was performed using Albacore v2.3.3 or Guppy v2.3.1 (Oxford Nanopore Technologies, Oxford, UK).

Genome sequence analyses were performed on a BioLinux v8 server (51) using publicly available genomic analysis tools wrapped together into a high-throughput genome analysis pipeline. Briefly, Illumina sequencing data were processed with Trim Galore v0.6.1 ([https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/](https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/)) to remove sequencing adaptors, low-quality bases, and poor-quality reads. Kraken v1 (52) taxonomic sequence classification of raw reads was used to confirm species designation, and to rule out

contamination. Illumina reads were assembled with SPAdes v3.11 to generate contigs with a 200 bp minimum length cutoff (53). Long-read sequence data generated for other studies (16, 23, 28, 41) were combined with Illumina data for the same isolate, and hybrid assembly was conducted using unicycler v0.4.7 or v0.4.8-beta (54). Assembled genomes were annotated using Prokka v1.14 and assembly quality was verified using QUAST (55). Genomes were included in the study if they had at least 35-fold Illumina read coverage, had assemblies with ≤ 350 contigs, and had total genome lengths ± 25% of the median of all isolates within each species group. Antimicrobial resistance and toxin genes were confirmed using BLASTn in line with EDS-HAT study phenotypic inclusion criteria. Specifically, all *S. aureus* genomes were confirmed to encode the *mecA* gene, all *E. faecalis* and *E. faecium* genomes were confirmed to encode a VanA or VanB operon, all *E. coli* and *Klebsiella* spp. genomes were confirmed to encode an identifiable extended-spectrum beta-lactamase (ESBL) enzyme, and all *C. difficile* genomes were confirmed to encode either toxin A and/or toxin B genes.

**Classification of genomospecies and lineages.** Within each species group, genome assemblies from this study and reference genome assemblies downloaded from the NCBI RefSeq database underwent pairwise average nucleotide identity (ANI) analysis using FastANI v1.3 (19). Genomes with ANI values >95% then underwent single-linkage hierarchical clustering using the hclust function from the R package stats v3.6. Each ANI cluster was manually assessed and assigned to a species based on the predominant nomenclature of genomes of type/reference strains within each cluster. Clusters that did not contain reference genomes, or where reference genomes were only named at the genus level, were named "genomospecies." Sequential numbers were appended to each uncharacterized genomospecies within each genus. Species identified using ANI and having greater than 100 isolates were further subdivided into clades and lineages based on multilocus sequence typing (ST), or phylogenetic analysis. STs were determined from assembled contigs using mlst v2 (https://github.com/tseemann/mlst). Species without a defined ST scheme (*P. mirabilis* and *S. marcescens*) were classified into clades or lineages by grouping isolates that shared <1000 core genome single nucleotide polymorphism (SNP) differences into the same lineage, with SNPs identified using snippy (https://github.com/tseemann/snippy). *Stenotrophomonas* genomospecies were named according to Gröschel et al. (56).

**Gene content and pangenome analyses.** Gene content matrices were obtained for all species groups with more than 50 isolates using the pangenome analysis program Roary v3.11 (57). Roary was run using a protein identity cutoff 80% for genera containing multiple species, and a cutoff 95% for individual species. Pangenome collector's curves were generated for each species group by calculating the number of unique genes present at increasing numbers of sampled genomes, with 1000 iterations of each sample size up to 250. Genetic clustering of genomes within species groups based on variable gene content was calculated and visualized using principal-component analysis of accessory genes (PCA-A) using the R packages vegan v2.5-7, and ggbiplot v0.55, with matrices of gene presence/absence used as input. Genes that were present in all isolates, present in only one isolate, or absent in only one isolate, were removed from analysis. PCA-A coordinate plots were visualized using GraphPad Prism version 7.0c.

**Core genome SNP comparisons, phylogenetic trees, evolutionary rate, and recombination analyses.** Within each genus, species, ST, or clade, SNPs were identified using snippy (https://github.com/tseemann/snippy). The most complete genome assembly (i.e., highest $N_{50}$) was used as a reference genome for SNP analysis. Core genome SNPs, defined as SNPs at nucleotide positions shared across all genomes in the sample group being compared, were used to calculate pairwise SNP distances and to generate maximum likelihood phylogenetic trees. Trees were generated with RAxML v8.2 using the general time reversible model of evolution (GTRCAT), Lewis correction for ascertainment bias, and 100 bootstrap replicates (58). Unless otherwise specified, reported SNP distances refer to core genome SNPs for all isolates belonging to the same ST. Pairwise SNP distances were visualized using the R package ggplot2 v3.3.5. Recombination and evolutionary rates were calculated for STs in four species groups (*P. aeruginosa*, *Clostridioides difficile*, VRE, and MRSA), and for STs within each group with more than 25 isolates. Estimates of relative recombination rates (R/Theta) and average size of recombinant sequences (delta) were assessed from core genome alignments using ClonalFrameML v1.12 (25), with default settings. The relative rate of recombination, which reflects the number of nucleotide changes introduced by recombination relative to each point mutation (r/m) was calculated as r/m = (R/Theta) × delta × $\nu$ (25), where $\nu$ is the average distance between recombined sequences. A core genome alignment and recombination-corrected phylogenetic tree were used to estimate evolutionary rates using TreeTime (24). Isolates that were found to be highly divergent from other isolates of the same ST (as revealed by an excess number of SNPs separating them from other isolates) were removed from the analysis.

**Antibiotic resistance gene detection and analysis.** Acquired antimicrobial resistance genes were detected by querying genome assemblies against the ResFinder database using BLASTn (26). A gene was considered present if the BLASTn percent identity multiplied by the sequence coverage was >80%. Resistance gene presence was mapped to a global phylogenetic tree constructed from amino acid sequences of 120 ubiquitous protein coding genes from the Genome Taxonomy Database Tool kit (59). Resistance gene co-occurrence was calculated using the %*% operator in R. This operator works by identifying the cross-products between any two genes found in a matrix of resistance genes identified in all isolates. The results were used to construct a relative frequency plot using the R package ggplot2 v3.3.5. To include only the most frequently co-occurring gene pairs in the plot, a relative frequency of 80% and a combined frequency of 50% were used as cutoff thresholds. Additionally, genes found in >250 isolates were excluded as they were suspected of not being acquired resistance genes. ESBL and carbapenemase enzyme distributions were determined by assigning enzyme types based on protein sequence, and only 100% protein sequence matches are reported.

**Shared sequence detection and analysis.** Putative mobile genetic elements were identified by searching for sequences >10kb that were present at high identity (>99.9%) in the genomes of isolates belonging to different species (<95% ANI) using nucmer (60). Sequences were organized into clusters using all-by-all BLASTn v2.7.1 (61), and clusters were visualized with Cytoscape v3.8.2 (62). Clustered shared sequences were determined as resembling plasmids, insertion sequences (ISs), transposons, prophages, or integrative conjugative elements by BLAST against complete plasmids from NCBI databases (63), MobileElementFinder (64), PHASTER (65), ProphET (66), and ICEberg (67), as well as comparison to the NCBI nr database and manual curation. Antimicrobial resistance genes in clustered sequences were identified by BLASTn against the ResFinder database (26). Clusters of orthologous groups of proteins (COG) categories were assigned to genes present in one or more clustered sequences, and the distribution of genes in each COG category was visualized with the pie function in R.

**Data availability.** Raw sequencing reads and genome assemblies were submitted to the NCBI Sequence Read Archive (SRA) and GenBank, with accession numbers listed in Table S2.

## SUPPLEMENTAL MATERIAL

Supplemental material is available online only.
**DATA SET S1**, XLSX file, 3.1 MB.
**DATA SET S2**, XLSX file, 1.6 MB.
**FIG S1**, PDF file, 0.8 MB.
**FIG S2**, PDF file, 0.05 MB.
**FIG S3**, PDF file, 0.9 MB.
**FIG S4**, PDF file, 0.04 MB.
**FIG S5**, PDF file, 1.5 MB.
**TABLE S1**, DOCX file, 0.02 MB.
**TABLE S2**, XLSX file, 0.3 MB.
**TABLE S3**, XLSX file, 0.01 MB.

## REFERENCES

1. Magill SS, O'Leary E, Janelle SJ, Thompson DL, Dumyati G, Nadle J, Wilson LE, Kainer MA, Lynfield R, Greissman S, Ray SM, Beldavs Z, Gross C, Bamberg W, Sievers M, Concannon C, Buhr N, Warnke L, Maloney M, Ocampo V, Brooks J, Oyewumi T, Sharmin S, Richards K, Rainbow J, Samper M, Hancock EB, Leaptrot D, Scalise E, Badrun F, Phelps R, Edwards JR, Emerging Infections Program Hospital Prevalence Survey T. 2018. Changes in prevalence of health care-associated infections in U.S. Hospitals. N Engl J Med 379:1732–1744. https://doi.org/10.1056/NEJMoa1801550.

2. Stone PW. 2009. Economic burden of healthcare-associated infections: an American perspective. Expert Rev Pharmacoecon Outcomes Res 9: 417–422. https://doi.org/10.1586/erp.09.53.

3. Centers for Disease Control and Prevention. 2020. Dec 02. Current HAI progress report. https://www.cdc.gov/hai/data/portal/progress-report.html. Accessed 05/05/2021.

4. Rice LB. 2008. Federal funding for the study of antimicrobial resistance in nosocomial pathogens: no ESKAPE. J Infect Dis 197:1079–1081. https://doi.org/10.1086/533452.

5. Centers for Disease Control and Prevention. 2013. ANTIBIOTIC RESISTANCE THREATS in the United States, 2013. CDC.

6. Lax S, Sangwan N, Smith D, Larsen P, Handley KM, Richardson M, Guyton K, Krezalek M, Shogan BD, Defazio J, Flemming I, Shakhsheer B, Weber S, Landon E, Garcia-Houchins S, Siegel J, Alverdy J, Knight R, Stephens B, Gilbert JA. 2017. Bacterial colonization and succession in a newly opened hospital. Sci Transl Med 9. https://doi.org/10.1126/scitranslmed.aah6500.

7. Curry SR, Muto CA, Schlackman JL, Pasculle AW, Shutt KA, Marsh JW, Harrison LH. 2013. Use of multilocus variable number of tandem repeats analysis genotyping to determine the role of asymptomatic carriers in Clostridium difficile transmission. Clin Infect Dis 57:1094–1102. https://doi.org/10.1093/cid/cit475.

8. Kanamori H, Rutala WA, Weber DJ. 2017. The role of patient care items as a fomite in healthcare-associated outbreaks and infection prevention. Clin Infect Dis 65:1412–1419. https://doi.org/10.1093/cid/cix462.

9. Santajit S, Indrawattana N. 2016. Mechanisms of antimicrobial resistance in ESKAPE pathogens. Biomed Res Int 2016:2475067. https://doi.org/10.1155/2016/2475067.

10. Quainoo S, Coolen JPM, van Hijum S, Huynen MA, Melchers WJG, van Schaik W, Wertheim HFL. 2017. Whole-genome sequencing of bacterial pathogens: the future of nosocomial outbreak analysis. Clin Microbiol Rev 30:1015–1063. https://doi.org/10.1128/CMR.00016-17.

11. Peacock SJ, Parkhill J, Brown NM. 2018. Changing the paradigm for hospital outbreak detection by leading with genomic surveillance of nosocomial pathogens. Microbiology (Reading) 164:1213–1219. https://doi.org/10.1099/mic.0.000700.

12. Sundermann AJ, Miller JK, Marsh JW, Saul MI, Shutt KA, Pacey M, Mustapha MM, Ayres A, Pasculle AW, Chen J, Snyder GM, Dubrawski AW,

Harrison LH. 2019. Automated data mining of the electronic health record for investigation of healthcare-associated outbreaks. Infect Control Hosp Epidemiol 40:314–319. https://doi.org/10.1017/ice.2018.343.

13. Miller JK, Chen J, Sundermann A, Marsh JW, Saul MI, Shutt KA, Pacey M, Mustapha MM, Harrison LH, Dubrawski A. 2019. Statistical outbreak detection by joining medical records and pathogen similarity. J Biomed Inform 91:103126. https://doi.org/10.1016/j.jbi.2019.103126.

14. Sundermann AJ, Chen J, Miller JK, Saul MI, Shutt KA, Griffith MP, Mustapha MM, Ezeonwuka C, Waggle K, Srinivasa V, Kumar P, Pasculle AW, Ayres AM, Snyder GM, Cooper VS, Van Tyne D, Marsh JW, Dubrawski AW, Harrison LH. 2020. Outbreak of Pseudomonas aeruginosa infections from a contaminated gastroscope detected by whole genome sequencing surveillance. Clin Infect Dis https://doi.org/10.1093/cid/ciaa1887.

15. Sundermann AJ, Babiker A, Marsh JW, Shutt KA, Mustapha MM, Pasculle AW, Ezeonwuka C, Saul MI, Pacey MP, Van Tyne D, Ayres AM, Cooper VS, Snyder GM, Harrison LH. 2020. Outbreak of vancomycin-resistant Enterococcus faecium in interventional radiology: detection through whole-genome sequencing-based surveillance. Clin Infect Dis 70:2336–2343. https://doi.org/10.1093/cid/ciz666.

16. Sundermann AJ, Chen J, Kumar P, Ayres AM, Cho ST, Ezeonwuka C, Griffith MP, Miller JK, Mustapha MM, Pasculle AW, Saul MI, Shutt KA, Srinivasa V, Waggle K, Snyder DJ, Cooper VS, Van Tyne D, Snyder GM, Marsh JW, Dubrawski A, Roberts MS, Harrison LH. 2021. Whole genome sequencing surveillance and machine learning of the electronic health record for enhanced healthcare outbreak detection. Clin Infect Dis https://doi.org/10.1093/cid/ciab946.

17. Roach DJ, Burton JN, Lee C, Stackhouse B, Butler-Wu SM, Cookson BT, Shendure J, Salipante SJ. 2015. A year of infection in the intensive care unit: prospective whole genome sequencing of bacterial clinical isolates reveals cryptic transmissions and novel microbiota. PLoS Genet 11: e1005413. https://doi.org/10.1371/journal.pgen.1005413.

18. Mosquera-Rendon J, Rada-Bravo AM, Cardenas-Brito S, Corredor M, Restrepo-Pineda E, Benitez-Paez A. 2016. Pangenome-wide and molecular evolution analyses of the Pseudomonas aeruginosa species. BMC Genomics 17:45. https://doi.org/10.1186/s12864-016-2364-4.

19. Jain C, Rodriguez RL, Phillippy AM, Konstantinidis KT, Aluru S. 2018. High throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. Nat Commun 9:5114. https://doi.org/10.1038/s41467-018-07641-9.

20. Jeukens J, Freschi L, Vincent AT, Emond-Rheault JG, Kukavica-Ibrulj I, Charette SJ, Levesque RC. 2017. A pan-genomic approach to understand the basis of host adaptation in achromobacter. Genome Biol Evol 9: 1030–1046. https://doi.org/10.1093/gbe/evx061.

21. Potter RF, Burnham CD, Dantas G. 2019. In silico analysis of Gardnerella genomospecies detected in the setting of bacterial vaginosis. Clin Chem 65:1375–1387. https://doi.org/10.1373/clinchem.2019.305474.

22. Roy PH, Tetu SG, Larouche A, Elbourne L, Tremblay S, Ren Q, Dodson R, Harkins D, Shay R, Watkins K, Mahamoud Y, Paulsen IT. 2010. Complete genome sequence of the multiresistant taxonomic outlier Pseudomonas aeruginosa PA7. PLoS One 5:e8842. https://doi.org/10.1371/journal.pone.0008842.

23. Marsh JW, Mustapha MM, Griffith MP, Evans DR, Ezeonwuka C, Pasculle AW, Shutt KA, Sundermann A, Ayres AM, Shields RK, Babiker A, Cooper VS, Van Tyne D, Harrison LH. 2019. Evolution of outbreak-causing carbapenem-resistant Klebsiella pneumoniae ST258 at a tertiary care hospital over 8 years. mBio 10. https://doi.org/10.1128/mBio.01945-19.

24. Sagulenko P, Puller V, Neher RA. 2018. TreeTime: maximum-likelihood phylodynamic analysis. Virus Evol 4:vex042. https://doi.org/10.1093/ve/vex042.

25. Didelot X, Wilson DJ. 2015. ClonalFrameML: efficient inference of recombination in whole bacterial genomes. PLoS Comput Biol 11:e1004041. https://doi.org/10.1371/journal.pcbi.1004041.

26. Zankari E, Hasman H, Cosentino S, Vestergaard M, Rasmussen S, Lund O, Aarestrup FM, Larsen MV. 2012. Identification of acquired antimicrobial resistance genes. J Antimicrob Chemother 67:2640–2644. https://doi.org/10.1093/jac/dks261.

27. Lerminiaux NA, Cameron ADS. 2019. Horizontal transfer of antibiotic resistance genes in clinical environments. Can J Microbiol 65:34–44. https://doi.org/10.1139/cjm-2018-0275.

28. Evans DR, Griffith MP, Sundermann AJ, Shutt KA, Saul MI, Mustapha MM, Marsh JW, Cooper VS, Harrison LH, Van Tyne D. 2020. Systematic detection of horizontal gene transfer across genera among multidrug-resistant bacteria in a single hospital. Elife 9. https://doi.org/10.7554/eLife.53886.

29. McDonnell G, Russell AD. 1999. Antiseptics and disinfectants: activity, action, and resistance. Clin Microbiol Rev 12:147–179. https://doi.org/10.1128/CMR.12.1.147.

30. McNally A, Oren Y, Kelly D, Pascoe B, Dunn S, Sreecharan T, Vehkala M, Valimaki N, Prentice MB, Ashour A, Avram O, Pupko T, Dobrindt U, Literak I, Guenther S, Schaufler K, Wieler LH, Zhiyong Z, Sheppard SK, McInerney JO, Corander J. 2016. Combined analysis of variation in core, accessory and regulatory genome regions provides a super-resolution view into the evolution of bacterial populations. PLoS Genet 12:e1006280. https://doi.org/10.1371/journal.pgen.1006280.

31. Inglin RC, Meile L, Stevens MJA. 2018. Clustering of pan- and core-genome of Lactobacillus provides novel evolutionary insights for differentiation. BMC Genomics 19:284. https://doi.org/10.1186/s12864-018-4601-5.

32. Freschi L, Vincent AT, Jeukens J, Emond-Rheault JG, Kukavica-Ibrulj I, Dupont MJ, Charette SJ, Boyle B, Levesque RC. 2019. The Pseudomonas aeruginosa pan-genome provides new insights on its population structure, horizontal gene transfer, and pathogenicity. Genome Biol Evol 11: 109–120. https://doi.org/10.1093/gbe/evy259.

33. Moradigaravand D, Boinett CJ, Martin V, Peacock SJ, Parkhill J. 2016. Recent independent emergence of multiple multidrug-resistant Serratia marcescens clones within the United Kingdom and Ireland. Genome Res 26:1101–1109. https://doi.org/10.1101/gr.205245.116.

34. Abreo E, Altier N. 2019. Pangenome of Serratia marcescens strains from nosocomial and environmental origins reveals different populations and the links between them. Sci Rep 9:46. https://doi.org/10.1038/s41598-018-37118-0.

35. Matteoli FP, Pedrosa-Silva F, Dutra-Silva L, Giachini AJ. 2021. The global population structure and beta-lactamase repertoire of the opportunistic pathogen Serratia marcescens. Genomics 113:3523–3532. https://doi.org/10.1016/j.ygeno.2021.08.009.

36. Eyre DW, Cule ML, Wilson DJ, Griffiths D, Vaughan A, O'Connor L, Ip CLC, Golubchik T, Batty EM, Finney JM, Wyllie DH, Didelot X, Piazza P, Bowden R, Dingle KE, Harding RM, Crook DW, Wilcox MH, Peto TEA, Walker AS. 2013. Diverse sources of C. difficile infection identified on whole-genome sequencing. N Engl J Med 369:1195–1205. https://doi.org/10.1056/NEJMoa1216064.

37. Coll F, Raven KE, Knight GM, Blane B, Harrison EM, Leek D, Enoch DA, Brown NM, Parkhill J, Peacock SJ. 2020. Definition of a genetic relatedness cutoff to exclude recent transmission of meticillin-resistant Staphylococcus aureus: a genomic epidemiology analysis. Lancet Microbe 1: e328–e335. https://doi.org/10.1016/S2666-5247(20)30149-X.

38. Miyoshi-Akiyama T, Tada T, Ohmagari N, Viet Hung N, Tharavichitkul P, Pokhrel BM, Gniadkowski M, Shimojima M, Kirikae T. 2017. Emergence and spread of epidemic multidrug-resistant Pseudomonas aeruginosa. Genome Biol Evol 9:3238–3245. https://doi.org/10.1093/gbe/evx243.

39. Didelot X, Eyre DW, Cule M, Ip CL, Ansari MA, Griffiths D, Vaughan A, O'Connor L, Golubchik T, Batty EM, Piazza P, Wilson DJ, Bowden R, Donnelly PJ, Dingle KE, Wilcox M, Walker AS, Crook DW, Peto TE, Harding RM. 2012. Microevolutionary analysis of Clostridium difficile genomes to investigate transmission. Genome Biol 13:R118. https://doi.org/10.1186/gb-2012-13-12-r118.

40. van Duin D, Arias CA, Komarow L, Chen L, Hanson BM, Weston G, Cober E, Garner OB, Jacob JT, Satlin MJ, Fries BC, Garcia-Diaz J, Doi Y, Dhar S, Kaye KS, Earley M, Hujer AM, Hujer KM, Domitrovic TN, Shropshire WC, Dinh A, Manca C, Luterbach CL, Wang M, Paterson DL, Banerjee R, Patel R, Evans S, Hill C, Arias R, Chambers HF, Fowler VG, Kreiswirth BN, Bonomo RA, Multi-Drug Resistant Organism Network I. 2020. Molecular and clinical epidemiology of carbapenem-resistant Enterobacterales in the USA (CRACKLE-2): a prospective cohort study. Lancet Infect Dis 20:731–741. https://doi.org/10.1016/S1473-3099(19)30755-8.

41. Babiker A, Evans DR, Griffith MP, McElheny CL, Hassan M, Clarke LG, Mettus RT, Harrison LH, Doi Y, Shields RK, Van Tyne D. 2020. Clinical and genomic epidemiology of carbapenem-nonsusceptible Citrobacter spp. at a tertiary health care center over 2 decades. J Clin Microbiol 58. https://doi.org/10.1128/JCM.00275-20.

42. World Health Organization. 22 September 2020 GLASS whole-genome sequencing for surveillance of antimicrobial resistance.

43. Ramirez MS, Tolmasky ME. 2010. Aminoglycoside modifying enzymes. Drug Resist Updat 13:151–171. https://doi.org/10.1016/j.drup.2010.08.003.

44. De Oliveira DMP, Forde BM, Kidd TJ, Harris PNA, Schembri MA, Beatson SA, Paterson DL, Walker MJ. 2020. Antimicrobial resistance in ESKAPE pathogens. Clin Microbiol Rev 33. https://doi.org/10.1128/CMR.00181-19.

45. Pal C, Bengtsson-Palme J, Kristiansson E, Larsson DG. 2015. Co-occurrence of resistance genes to antibiotics, biocides and metals reveals novel

insights into their co-selection potential. BMC Genomics 16:964. https://doi.org/10.1186/s12864-015-2153-5.

46. Canton R, Gonzalez-Alba JM, Galan JC. 2012. CTX-M enzymes: origin and diffusion. Front Microbiol 3:110.

47. Redondo-Salvo S, Fernandez-Lopez R, Ruiz R, Vielva L, de Toro M, Rocha EPC, Garcillan-Barcia MP, de la Cruz F. 2020. Pathways for horizontal gene transfer in bacteria revealed by a global map of their plasmids. Nat Commun 11:3602. https://doi.org/10.1038/s41467-020-17278-2.

48. Partridge SR, Kwong SM, Firth N, Jensen SO. 2018. Mobile genetic elements associated with antimicrobial resistance. Clin Microbiol Rev 31. https://doi.org/10.1128/CMR.00088-17.

49. Chandrangsu P, Rensing C, Helmann JD. 2017. Metal homeostasis and resistance in bacteria. Nat Rev Microbiol 15:338–350. https://doi.org/10.1038/nrmicro.2017.15.

50. Villapun VM, Dover LG, Cross A, Gonzalez S. 2016. Antibacterial metallic touch surfaces. Materials (Basel) 9:736. https://doi.org/10.3390/ma9090736.

51. Field D, Tiwari B, Booth T, Houten S, Swan D, Bertrand N, Thurston M. 2006. Open software for biologists: from famine to feast. Nat Biotechnol 24:801–803. https://doi.org/10.1038/nbt0706-801.

52. Wood DE, Salzberg SL. 2014. Kraken: ultrafast metagenomic sequence classification using exact alignments. Genome Biol 15:R46. https://doi.org/10.1186/gb-2014-15-3-r46.

53. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Prjibelski AD, Pyshkin AV, Sirotkin AV, Vyahhi N, Tesler G, Alekseyev MA, Pevzner PA. 2012. SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. J Comput Biol 19:455–477. https://doi.org/10.1089/cmb.2012.0021.

54. Wick RR, Judd LM, Gorrie CL, Holt KE. 2017. Unicycler: resolving bacterial genome assemblies from short and long sequencing reads. PLoS Comput Biol 13:e1005595. https://doi.org/10.1371/journal.pcbi.1005595.

55. Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUAST: quality assessment tool for genome assemblies. Bioinformatics 29:1072–1075. https://doi.org/10.1093/bioinformatics/btt086.

56. Groschel MI, Meehan CJ, Barilar I, Diricks M, Gonzaga A, Steglich M, Conchillo-Sole O, Scherer IC, Mamat U, Luz CF, De Bruyne K, Utpatel C, Yero D, Gibert I, Daura X, Kampmeier S, Rahman NA, Kresken M, van der Werf TS, Alio I, Streit WR, Zhou K, Schwartz T, Rossen JWA, Farhat MR, Schaible UE, Nubel U, Rupp J, Steinmann J, Niemann S, Kohl TA. 2020. The phylogenetic landscape and nosocomial spread of the multidrug-resistant opportunist Stenotrophomonas maltophilia. Nat Commun 11:2044. https://doi.org/10.1038/s41467-020-15123-0.

57. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, Fookes M, Falush D, Keane JA, Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis. Bioinformatics 31:3691–3693. https://doi.org/10.1093/bioinformatics/btv421.

58. Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics 30:1312–1313. https://doi.org/10.1093/bioinformatics/btu033.

59. Parks DH, Chuvochina M, Waite DW, Rinke C, Skarshewski A, Chaumeil PA, Hugenholtz P. 2018. A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. Nat Biotechnol 36:996–1004. https://doi.org/10.1038/nbt.4229.

60. Marcais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. 2018. MUMmer4: a fast and versatile genome alignment system. PLoS Comput Biol 14:e1005944. https://doi.org/10.1371/journal.pcbi.1005944.

61. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. 1990. Basic local alignment search tool. J Mol Biol 215:403–410. https://doi.org/10.1016/S0022-2836(05)80360-2.

62. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T. 2003. Cytoscape: a software environment for integrated models of biomolecular interaction networks. Genome Res 13:2498–2504. https://doi.org/10.1101/gr.1239303.

63. Che Y, Yang Y, Xu X, Brinda K, Polz MF, Hanage WP, Zhang T. 2021. Conjugative plasmids interact with insertion sequences to shape the horizontal transfer of antimicrobial resistance genes. Proc Natl Acad Sci U S A 118.

64. Johansson MHK, Bortolaia V, Tansirichaiya S, Aarestrup FM, Roberts AP, Petersen TN. 2021. Detection of mobile genetic elements associated with antibiotic resistance in Salmonella enterica using a newly developed web tool: MobileElementFinder. J Antimicrob Chemother 76:101–109. https://doi.org/10.1093/jac/dkaa390.

65. 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. https://doi.org/10.1093/nar/gkw387.

66. Reis-Cunha JL, Bartholomeu DC, Manson AL, Earl AM, Cerqueira GC. 2019. ProphET, prophage estimation tool: a stand-alone prophage sequence prediction tool with self-updating reference database. PLoS One 14:e0223364. https://doi.org/10.1371/journal.pone.0223364.

67. Bi D, Xu Z, Harrison EM, Tai C, Wei Y, He X, Jia S, Deng Z, Rajakumar K, Ou HY. 2012. ICEberg: a web-based resource for integrative and conjugative elements found in Bacteria. Nucleic Acids Res 40:D621–6. https://doi.org/10.1093/nar/gkr846.