

Title: Phylogenomic study of *Staphylococcus aureus* and *Staphylococcus haemolyticus* clinical isolates from Egypt

Authors: Cesar Montelongo^a, Carine R. Mores^a, Catherine Putonti^{a,b,c}, Alan J. Wolfe^a, Alaa Abouelfetouh^{d,e*}

Affiliations:

^aDepartment of Microbiology and Immunology, Stritch School of Medicine, Loyola University Chicago, Maywood, IL 60153 USA

^bBioinformatics Program, Loyola University Chicago, Chicago, IL 60660 USA

^cDepartment of Biology, Loyola University Chicago, Chicago, IL 60660 USA

^dDepartment of Microbiology and Immunology, Faculty of Pharmacy, Alexandria University, Alexandria, Egypt

^eDepartment of Microbiology and Immunology, Faculty of Pharmacy, Alalamein International University, Alamein, Egypt

Author email addresses:

Cesar Montelongo:

cmontelongoherna@luc.edu

Carine R. Mores:

carinermores@gmail.com

Catherine Putonti:

cputonti@luc.edu

Alan J. Wolfe:

Awolfe@luc.edu

24 Alaa Abouelfetouh:

25 Alaa.abouelfetouh@pharmacy.alexu.edu.eg

26 aabouelfetouh@Aiu.edu.eg

27

28 **Corresponding Author:** Alaa Abouelfetouh; Alaa.abouelfetouh@pharmacy.alexu.edu.eg

29

30 **Keywords:** *Staphylococcus aureus*; *Staphylococcus haemolyticus*; Middle East; MLST

31

32 **Repositories:** Raw sequencing reads and assembled genomes can be found at BioProject

33 Accession number PRJNA648411 (<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA648411>).

34

Abstract

Antibiotic resistant *Staphylococcus* infections are a global concern, with increasing cases of resistant *Staphylococcus aureus* and *Staphylococcus haemolyticus* found circulating in the Middle East. While extensive surveys have described the prevalence of resistant infections in Europe, Asia, and North America, the population structure of resistant staphylococcal Middle Eastern clinical isolates is poorly characterized. We performed whole genome sequencing of 56 *S. aureus* and 10 *S. haemolyticus* isolates from Alexandria Main University Hospital. Supplemented with additional publicly available genomes from the region (34 *S. aureus* and 6 *S. haemolyticus*), we present the largest genomic study of staphylococcal Middle Eastern isolates. These genomes include 20 *S. aureus* multilocus sequence typing (MLST) types and 9 *S. haemolyticus* MLSTs, including 3 and 1 new MLSTs, respectively. Phylogenomic analyses of each species core genome largely mirrored MLSTs, irrespective of geographical origin. The hospital-acquired *spa* t037/SCCmec III/MLST CC8 clone represented the largest clade, comprising 22% of *S. aureus* isolates. Similar to other regional genome surveys of *S. aureus*, the Middle Eastern isolates have an open pangenome, a strong indicator of gene exchange of virulence factors and antibiotic resistance genes with other reservoirs. We recommend stricter implementation of antibiotic stewardship and infection control plans in the region.

Impact Statement

Staphylococci are under-studied despite their prevalence within the Middle East. Methicillin-resistant *Staphylococcus aureus* (MRSA) is endemic to hospitals in this region, as are other antibiotic-resistant strains of *S. aureus* and *S. haemolyticus*. To provide insight into the strains currently in circulation within Egypt, we performed whole genome sequencing of 56 *S. aureus* and 10 *S. haemolyticus* isolates from Alexandria Main University Hospital (AMUH). Through analysis of these genomes, as well as other genomes of isolates from the Middle East, we were able to produce a more complete picture of the current diversity than traditional molecular typing strategies. Furthermore, the *S. haemolyticus* genome analyses provide the first insight into strains found in Egypt. Our analysis of resistance and virulence mechanisms carried by these strains provides invaluable insight into future plans of antibiotic stewardship and infection control within the region.

Data Summary

Raw sequencing reads and assembled genomes can be found at BioProject Accession number PRJNA648411 (<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA648411>).

Introduction

Staphylococci are a heterogenous group of commensal bacteria in humans with the potential to cause infections ¹. Two staphylococcal species especially relevant to the clinical setting are *Staphylococcus aureus* and *Staphylococcus haemolyticus*. *S. aureus* is arguably the most clinically important staphylococcal species; the infections it can cause range from mild erythema to serious life-threatening ailments, including septicemia, pneumonia, and endocarditis ². A difficulty in treating and controlling *S. aureus* stems from its prevalence and increasing resistance to clinically used antibiotics, resulting in it being one of the leading agents for nosocomial and community-acquired infections ^{3,4}. *S. haemolyticus* is the second most common staphylococcal species isolated in human blood cultures and a prominent reservoir for antibiotic resistance genes, which can be shared with other Staphylococci, including *S. aureus* ⁵⁻⁷.

Epidemiological surveillance and profiling are key to managing Staphylococci ^{8,9}. Historically, profiling of Staphylococci has relied on complementary molecular typing strategies, such as Multi Locus Sequence Typing (MLST), typing of hypervariable short repeats in Protein A (*spa*), subtyping elements in the cassette chromosome *mec* (*SCCmec*), and presence of the Panton-Valentine leucocidin (PVL) ¹⁰. MLST consists of comparing the sequence of specific housekeeping genes in bacteria; the strategy is effective at tracking a broad range of clones over a global area, but prior to whole-genome sequencing (WGS) this method was slow and expensive ¹⁰. *spa* typing complements MLST by tracking the molecular evolution of *S. aureus*, given the relevance of Protein A to the infectious process ⁹. *SCCmec* permits profiling of clinically relevant antibiotic resistances, including *mecA*, which results in methicillin-resistant *S. aureus* (MRSA) ^{11,12}. PVL is an important cytotoxin in MRSA that is common in community-acquired MRSA (CA-MRSA), but uncommon in hospital-associated MRSA (HA-MRSA) ¹¹. Traditionally, these profiling strategies have been a powerful means to type, trace, and manage

Staphylococcal infections, but technical limitations curtail the usefulness of molecular typing in real time^{9,10}.

The increasing utility, speed, and inexpensiveness of WGS in the clinical setting is poised to immensely benefit Staphylococci profiling^{10,13–15}. WGS allows access to the entire staphylococcal genome, including sequence data for typing MLST, *spa*, *SCCmec*, and PVL. In addition, WGS allows us to study the phylogenomic lineage, core, and accessory genome of isolates from an infectious outbreak or a geographical area. A key question is how WGS analysis compares to traditional typing techniques. For example, there is evidence that phylogenomic data does not always agree with standard typing methods; skepticism also exists that WGS can reliably detect single nucleotide polymorphisms (SNPs) in sensitive genetic content¹⁰. In contrast, studies have demonstrated that WGS can be used to type, discriminate, and cluster staphylococcal isolates for the purpose of outbreak control^{13–15}. WGS could be used to close the gap in staphylococcal management in regions that have not been extensively monitored, such as the Middle East and specifically Egypt.

The epidemiology of Staphylococci in non-European countries of the Mediterranean region is under-studied¹⁶. Antibiotic resistance in *S. haemolyticus* has been identified in Middle Eastern countries, such as Turkey⁶ and Egypt⁷. There is evidence that MRSA is prevalent and endemic to hospitals in this region, with a median MRSA prevalence of 38% in Algeria, Cyprus, Egypt, Jordan, Lebanon, Malta, Morocco, Tunisia and Turkey¹⁷. Broadly speaking, PVL prevalence is reported as low in some of these countries, indicating a predominance of HA-MRSA^{16,18,19}. Research into the lineage of Staphylococci in this region is urgent, as it would give us both a present and future assessment of staphylococcal epidemiology.

Generally, molecular typing and phylogeny data are limited from this region. Multiple isolates in Palestine were typed as ST22 with a minority typed as ST80-MRSA-IV and PVL-positive²⁰. In Jordan, genotyping of *S. aureus* isolates revealed that the majority were ST80-MRSA-IV²¹. In Lebanon, the primary lineage was PVL-positive ST80-MRSA-IV followed by PVL-positive ST30-MSSA¹⁷. In Algeria, it was reported that ST80-MRSA-IV was present in most neonates tested over an 18-month period, with a minority of these PVL-positive²². Finally, for Egypt, it has been reported that the prevalent MLSTs are ST30, ST80, and a novel type, ST1010; PVL prevalence has been estimated at 19%²³. Enany *et al.* reported that the Egyptian ST80 lineage was different from the globally prevalent ST80, primarily due to a unique *spa* type and antimicrobial resistance²³.

Egypt presents a unique case-study for staphylococcal distribution in Arab countries²⁴. Egypt's cultural and geographical placement may facilitate local Staphylococcal exposure to international lineages, both from the Middle East and elsewhere. The accessibility of WGS presents an opportunity to profile Staphylococci in Egypt and the rest of the Arab region in terms of gene marker typing, core genome, and phylogenomics. Prior to this study, there were limited genomic data of *S. aureus* and *S. haemolyticus* in this region.

Here, we report the phylogenetic and phylogenomic associations of 56 *S. aureus* and 10 *S. haemolyticus* isolates from Egypt and their relationship to 34 *S. aureus* and 6 *S. haemolyticus* isolates from Egypt, Kuwait, Lebanon, Tunisia, Palestine, United Arab Emirates, Morocco, and Sudan. WGS afforded insight into the lineage and genetic content of these two staphylococcal species, including type information historically obtained using molecular methods. Both the MLST and SCCmec type mirrored the core genome, indicating that WGS is a fast and accessible option for Staphylococcal profiling. We identified multiple MLST and clonal complexes in

circulation in the region, including 3 new genotypes. Genome analysis indicated that *S. aureus* in Egypt has an open pangenome that includes virulence genes in both the core and accessory genomes. Surveillance and profiling of Staphylococci are key to infection control, and we have shown that WGS can be a valuable asset, especially in regions where Staphylococci have not been well studied, such as the Middle East.

Results

S. aureus and *S. haemolyticus* isolates were collected from patients presenting to the Medical Microbiology Laboratory at AMUH between September and December 2015. Draft genomes for 66 of the clinical isolates were of high quality and were included in our analysis. These genomes included 56 *S. aureus* and 10 *S. haemolyticus* isolates, assembled on average into 71 and 126 contigs, respectively. Additional publicly available *S. aureus* and *S. haemolyticus* strains were identified and included in subsequent analyses: 34 *S. aureus* (from Egypt n=17; Kuwait n=5; Lebanon n=4; Tunisia, Palestine and United Arab Emirates n=2 each; Morocco and Sudan n=1 each) and 6 *S. haemolyticus* (all from Egypt). **Supplementary Table S1** lists the available metadata and presents the genome assembly statistics for *S. aureus* and *S. haemolyticus*. The *S. aureus* genomes were, on average, larger than *S. haemolyticus* genomes: 2.8 Mbp and 2.5 Mbp, respectively. Genome size and GC content were on par with other publicly available genomes. Each draft genome was annotated using NCBI's PGAP, identifying an average of 2,839 and 2,495 coding sequences (CDS) for *S. aureus* and *S. haemolyticus*, respectively. The strains varied in their number of rRNA operons and tRNAs.

Strain genotyping

The genomes represent varied MLSTs. The 16 *S. haemolyticus* isolates examined here belonged to nine MLSTs, including a new genotype ST-74 (strain 51) assigned as a result of this study, and an isolate of unknown ST (strain 7A). ST-3 was the most common amongst the isolates examined (n=4) (**Supplementary Table S2**). A total of 20 *S. aureus* MLSTs were identified, including three novel types ST-5860 (strain 48), ST-5861 (strain 2705404), and ST-5862 (strain 2705410); all three of these strains came from prior studies and were isolated from Egypt, Kuwait and Lebanon, respectively (**Supplementary Table S2**). Twelve different MLSTs were identified among the Egyptian isolates; ST-239 was the most prevalent (n=24), followed by ST-1 (n=19) and then ST-80 (n=12). Two of the AMUH *S. aureus* isolates, strains AA32 and AA35, could not be typed due to incomplete sequences.

S. aureus isolates could be categorized into seven clonal complexes (CC) (**Supplementary Table S2**), the largest being CC8 (n=26), consisting mainly of Egyptian isolates and one Moroccan isolate (strain 12480433). 12 strains were identified as ST-80, which does not belong to a clonal complex. 20 different *spa* types were identified in addition to 7 isolates that could not be typed. The predominant *spa* type was t037 (n=33), with all but one belonging to CC8; 30 of these 33 isolates belonged to SCCmec III. The next most frequent *spa* type was t127 (n=19), all belonging to CC1. **Table 1** summarizes these results.

Core and pangenomes of S. haemolyticus and S. aureus strains

To investigate the core genome and pangenome of *S. haemolyticus*, we added the 6 publicly available *S. haemolyticus* genomes from Egypt to our 10 *S. haemolyticus* genomes (**Supplementary Table S1**). The pangenome for these strains included 3,541 genes (**Supplementary Fig. S1**), with 1,834 single copy number genes in the core genome. Included

within these core genes are the virulence factors autolysin (*atl*), elastin binding protein (*ebp*),
thermonuclease (*nuc*), and cytolyisin (*cylR2*).

In addition to our 56 *S. aureus* genomes, our core and pangenome analysis included 34 *S. aureus* draft genome assemblies from the Arab region (**Supplementary Table S1**). The pangenome of these 90 isolates contained 4,283 genes (**Fig. 1, panel A**), the core genome included 1,501 single copy number genes, and the accessory genome contained 2,178 genes. These analyses show that the Arab isolates have an open pangenome. The functionality of the genes within the *S. aureus* core genome was determined according to their COG categories (**Fig. 1, panel B**). The core genome was further examined for virulence factors, finding the same gene related to autolysin (*atl*) that was found in the *S. haemolyticus* core. We also identified genes associated with intercellular adhesin, cysteine protease, thermonuclease, capsule, and the Type VII secretion system (**Table 2**).

In addition to the virulence factors found within the core genome, we identified virulence factors and antibiotic resistance genes within the accessory genome (**Supplementary Tables S3 and S4**). The isolates were screened for the presence of *lukF/S-PV*, which encodes PVL. Isolates positive for PVL were mainly (77%) *mecA* positive, present in CC1, ST-80, CC30 and CC8, and obtained from Egypt, Kuwait, Tunisia, Lebanon and Morocco. Importantly, isolates obtained from CA infections belonged to CC1, ST-80, CC5, CC97 and CC8, making PVL presence a good predictor for the ability of the isolate to cause CA infections.

Phylogenomic study of S. haemolyticus and S. aureus strains from the Arab region

The core genes were used to derive phylogenies for each species. The *S. haemolyticus* isolates were all from Egypt and clustered into two clades (**Fig. 2**). As the tree shows, variation between the core genomes of these isolates was minor. Furthermore, the clade structure of the genomes corresponded with MLST, indicated in the bar of **Fig. 2**. The MLST tree for these genomes is shown in **Supplementary Fig. S2**.

S. aureus isolates came from all over the region, and clustered into six clades, with Egyptian isolates represented in all clades (**Fig. 3 and 4**). Clade 1 isolates belonged to ST-1 and were from Egypt and the UAE, clade 2 contained the majority of the Arab isolates including some from Egypt. The predominating clone seen among 46.7% of the isolates within clade 2 was *spa* t044/*SCCmec* IV/ST-80, which shows some degree of shared content between these isolates. Clade 3 isolates were solely from Egypt and belonged mainly to ST-15 and ST-5. Clade 4 comprised isolates from Egypt, Sudan and Palestine, with the majority belonging to ST-22 and ST-361. Clade 5 contained isolates from Egypt, belonging mainly to ST-97. The remaining isolates were in clade 6, of ST-239 and from Egypt, with the exception of one Moroccan isolate. This clade represents a *spa* t037/*SCCmec* III/MLST CC8 clone. The phylogenetic tree derived from the core genome sequences corresponded with the tree derived from the MLST marker genes (**Fig. 3 and Supplementary Fig. S2**). 14 isolates lacked *mecA* (**Fig. 4**, pale green star) and occurred predominantly in CC1 (n=5), CC15 (n=3) and CC30, CC8 and ST-80, with one isolate in each; in addition, three isolates belonged to ST-361 (n=2) or ST-5860 (n=1). 13 of these *mecA* negative isolates were from Egypt.

Discussion

S. aureus is a major human pathogen in hospital and community settings, with the infection rate of MRSA increasing on a global scale at varying rates^{25,26}. While extensive surveys have provided insight into the prevalence of such resistant infections in Europe²⁷, Asia^{28,29}, and North America^{30,31}, limited data is available for the Arab region³². Furthermore, antibiotic resistant *S. haemolyticus* strains have been identified worldwide⁵, including Turkey⁶ and Egypt⁷. Prior to the study initiated here, there were limited genomic data for these two *Staphylococcus* species from the Middle East. The addition of 56 *S. aureus* and 10 *S. haemolyticus* genomes enabled our investigation of strain diversity within the region. With the majority (or all, in the case of *S. haemolyticus*) of genomes representing isolates from Egypt, we could investigate members of these species currently in circulation. We found that several different MLSTs and clonal complexes are in circulation within Egypt and more broadly within the region. 20 *S. aureus* MLSTs were identified in the region, including 3 new genotypes identified here, and 12 of these are in circulation within Egypt. Analysis of the *S. haemolyticus* genomes found 9 MLSTs in circulation within Egypt, including one new genotype.

The core genome for the *S. aureus* strains is slightly larger than that previously calculated for the species^{33,34}. This is expected, however, as our analysis is restricted to fewer genomes from a single region. Both our *S. haemolyticus* and *S. aureus* core genomes include the gene *atl*, which is essential for biofilm formation³⁵. Furthermore, the *ica* gene cluster, also associated with biofilm formation³⁶, as well as its regulator *icaR*³⁷, are included in the core genome of the *S. aureus* strains examined here. The presence of *atl* and the *ica* gene clusters signifies the biofilm potential of the isolates. This potential is relevant because 80% of human microbial

infections are complicated by biofilm formation, such as in wounds, IV catheters, sutures and implants (see reviews^{38,39}). Moreover, the biofilm capacity to evade the host's immune defenses, the inability of most antibiotic treatment regimens to eradicate existing biofilms and the fact that biofilms serve as a good medium for exchange of genetic material (e.g. plasmids) between cells make biofilm formation a major health concern in the clinical setting (⁴⁰; see reviews^{41,42}).

The Arab *S. aureus* genomes have an open pangenome, evidence of gene exchange between these isolates and other reservoirs. *S. aureus* is naturally competent⁴³, and horizontal gene transfer (HGT) between strains, coagulase-negative *Staphylococcus* (CoNS) strains, and other species is well documented (see review⁴⁴). Recently, HGT was shown to be a driver of persistent *S. aureus* infections within patients⁴⁵. Genes within the accessory genome included virulence factors and antibiotic resistance genes (**Supplementary Tables S3 and S4**). Prior comparative genomic studies for this species similarly found an open pangenome and resistance genes within the accessory genome³³.

Phylogenomic analyses of the core genome largely mirrored MLST types (**Fig. 3**). This was irrespective of geographical origin. Interestingly, strains of the same SCCmec type had a more similar core genome sequence (**Fig. 4**). In a prior phylogenetic study, John and co-workers found that 16S rRNA gene sequence similarity did not correspond with SCCmec type, leading them to conclude that horizontal gene transfer plays a role in resistance gene acquisition³³. However, only two SCCmec types were identified for the samples examined here. Recently, Soliman *et al.* published a study characterizing the genomes of 18 MRSA isolates from a tertiary care hospital in Cairo, Egypt; their isolates were primarily SCCmec types V (n=9) and VI (n=2), not observed within our larger collection⁴⁶. Similarly, SCCmec type V and IV have been most frequently observed in other *S. aureus* studies within the region⁴⁷⁻⁴⁹. Rather, our study found

that SCCmec type III and IV were equally prevalent within the region. SCCmec type III predominated among HA-MRSA infections, suggesting that, in contrast to these prior studies, our isolates indicate that HA-infections have a higher incidence among the patients tested here. AMUH, where our isolates were collected, is the largest tertiary hospital and main referral center in the Northern sector of Egypt; thus, patients with more severe infections would be more likely to be treated at AMUH than at any other hospital in the region. Prior studies have found SCCmec type III to be the predominant type in Asian countries²⁸. Besides, the SCCmec type III/MLST ST-239 is the oldest pandemic strain of MRSA⁵⁰, which might explain its prevalence among the current collection of isolates.

The *S. haemolyticus* and *S. aureus* genomes examined here provide insight into the diversity of strains currently in circulation in Egypt, particularly with respect to their encoded virulence factors and antibiotic resistance genes. WGS analysis enabled a more complete picture of this diversity than molecular typing strategies. The *S. haemolyticus* genomes provide the first insight into strains found in Egypt. Identifying the main genotypes, as well as the resistance and virulence mechanisms among the resistant isolates in the region, can drive antibiotic stewardship and infection control plans.

Methods

Bacterial isolates

A total of 89 *S. aureus* and 14 *S. haemolyticus* consecutive non-duplicate isolates were collected from the Medical Microbiology Laboratory at Alexandria Main University Hospital (AMUH) between September and December 2015. These isolates were obtained from various clinical specimens, including pus, blood, sputum, urine, tissue, aspirate and broncho-alveolar lavage

(BAL). The identity of the isolates was determined using conventional methods, such as Gram staining, growth on and fermentation of mannitol salt agar, growth on DNase agar and slide coagulase testing using Dryspot Staphytest Plus (Oxoid Ltd, England), and confirmed using Matrix Assisted Laser Desorption Ionization - Time of Flight Mass Spectrometry (MALDI-TOF MS) (Bruker Daltonik, USA). The isolates were further classified as either hospital-acquired or community-acquired infections based on a 48 h window between the dates of patient admission and isolate collection ⁵¹.

DNA extraction

Colonies grown on tryptone soya agar (TSA) plates were harvested and washed in 1 ml phosphate buffer saline (PBS) and resuspended in 0.5 ml SET (75mM NaCl, 25mM EDTA, 20mM Tris, pH 7.5), to which 50ul of fresh 20 mg/ml lysozyme in PBS and 30ul Mutanolysin were added; the mixture was incubated at 37°C for 60 min. The cells were then treated with 60ul 10% sodium dodecyl sulphate and 20ul proteinase K and incubated at 55°C for two hours with gentle inversion. The suspension was mixed gently with 210ul of 6M NaCl, and 700ul phenol:chloroform were added, followed by incubation at room temperature for 30-60 minutes, using a rotating wheel for gentle mixing. The suspension was then centrifuged at maximum speed for 10 min and the aqueous phase was transferred to a new microfuge tube and mixed gently with an equal volume of isopropanol. The tubes were centrifuged to produce a DNA pellet that was washed with 70% ethanol, which was left to evaporate overnight. The pellets were resuspended in 50ul ddH₂O and stored at -20°C till further processing.

Genome Sequencing and Genome Assembly

The Illumina Nextera kit was used for whole genome library preparation. Each isolate was sequenced using the Illumina MiSeq System, producing paired-end 2x250 bp reads. Quality control and de-multiplexing of sequence data was done with onboard MiSeq Control software and MiSeq Reporter v3.1. Raw reads were trimmed using Sickle v1.33 (<https://github.com/najoshi/sickle>) and assembled using SPAdes v3.13.0⁵² with the “only-assembler” option for k=55, 77, 99, and 127. Genome coverage was calculated using BBMap v38.47 (<https://sourceforge.net/projects/bbmap/>). Contigs shorter than 500 bp were pruned using bioawk (<https://github.com/lh3/bioawk>). Genome assemblies were annotated using PATRIC v3.3.18⁵³. Genomes were deposited in NCBI’s Assembly database, along with raw sequence data in SRA under BioProject PRJNA648411. Deposited genomes were annotated using the NCBI Prokaryotic Genome Annotation Pipeline (PGAP) v5.0⁵⁴. Unless previously noted, default parameters were used for each software tool. To complement our analysis of the genomes from AMUH, raw sequence data for 41 *S. aureus* and 10 *S. haemolyticus* strains were retrieved from NCBI. These records were identified by searching SRA (as of January 2020) for strains isolated in the Arab region. These raw reads were processed as indicated above. High-quality assemblies were included in subsequent analyses.

Bioinformatic Analysis

Multilocus sequence typing (MLST) was determined using the MLST v2.0.4 web server available through the Center for Genomic Epidemiology⁵⁵. MLST allele sequence and profile data were obtained from PubMLST v2.0.0⁵⁶. *spa* typing was performed using the online tool SpaTyper v1.0 available through the Center for Genomic Epidemiology⁵⁷. SCCmec typing was performed using SCCmecFinder v1.2 online tool available through the Center for Genomic

Epidemiology (<https://cge.cbs.dtu.dk/services/SCCmecFinder/>)^{58,59}. Resistance and virulence genes were identified using PATRIC v3.6.5⁶⁰ and VFAnalyzer⁶¹.

Phylogenomic and Phylogenetic Analysis

The core and pangenomes were generated using anvi'o v5.1. The following scripts were used to calculate the pangenome: anvi-gen-genomes-storage, anvi-pan-genome and anvi-display-pan, and the following script was used to calculate the core genome: anvi-get-sequences-for-gene-clusters^{62,63}. Functional groups for the core genome were determined by querying core genome amino acid sequences against the COG database⁶⁴ through anvi'o. The core genes were concatenated for each genome and then aligned using MAFFT v7.388⁶⁵. The tree was built using the FastTree v2⁶⁶ plugin in Geneious Prime v2019.2 (Biomatters Ltd., Auckland, New Zealand). MLST ST sequences were downloaded from PubMLST v2.0.0⁵⁶, aligned in Geneious Prime v2019.2 and the trees were built using the FastTree v2⁶⁶ plugin in Geneious Prime v2019.2. iTOL v5.6.1⁶⁷ was used to annotate and visualize all trees.

Author Statements

Acknowledgments

We acknowledge Roberto Limeira and Loyola Genomics Facility for performing the whole genome sequencing of the isolates. We also acknowledge funding from NIH (R01 DK104718 awarded to AJW), NSF (1661357 awarded to CP), USAID (GSP-T85 awarded to AA) and DFG (ZI 665/3-1 awarded to AA). The funders did not play a part in the design or conduct of the study

367 Authors and Contributors

368 CM: Formal Analysis; Writing – Original Draft Preparation; Writing – Review and Editing

369 CRM: Formal Analysis; Writing – Original Draft Preparation; Writing – Review and Editing

370 CP: Formal Analysis; Writing – Original Draft Preparation; Visualization; Writing – Review and
371 Editing

372 AJW: Formal Analysis; Conceptualisation; Writing – Review and Editing

373 AA: Conceptualisation; Formal Analysis; Writing – Original Draft Preparation; Writing –
374 Review and Editing

375 All authors reviewed the manuscript.

376 Competing Interests

377 AJW is a member of the Advisory Board of Urobiome Therapeutics. The remaining authors
378 report no disclosures.

379 Data availability

380 Raw sequencing reads and assembled genomes can be found at BioProject Accession number
381 PRJNA648411 (<https://www.ncbi.nlm.nih.gov/bioproject/PRJNA648411>)

382 Ethics declarations:

383 Ethics approval and consent to participate

384 Not applicable

385

386 Consent for publication

387 Not applicable

388

389 Figure legends

Figure 1. Genome analysis of 90 Arab *S. aureus* strains. **(a)** The pangenome. Each ring corresponds to a single genome. Each radial extension in the ring corresponds to the presence (black) or absence (light gray) of a given gene cluster (homologous gene). The bar charts list the number of genes identified in the given genome (top) and the number of singleton genes or genes that are unique to the given genome (bottom). The pangenome of these 90 isolates contained 4,283 genes, the core genome included 1,501 single copy number genes, and the accessory genome contained 2,178 genes. **(b)** Functionality of genes contained within the core genome. The same autolysin gene (*atl*) found in the core genome of *S. haemolyticus* was found in *S. aureus*.

Figure 2. Phylogeny based upon the core genes for the *S. haemolyticus* isolates. All *S. haemolyticus* isolates were from Egypt and clustered into two clades corresponding with MLST.

Figure 3. *S. aureus* core genome phylogeny colored by geographical origin of isolation (strain name color) and MLST (right bar). *S. aureus* isolates were from different parts of the region, and clustered into six clades, each containing Egyptian isolates. Clade 1 isolates belonged to ST-1 and were from Egypt and the UAE, clade 2 contained the majority of the Arab isolates, with *spa* t044/SCCmec IV/ST-80 as the predominating clone. Clade 3 isolates were solely from Egypt and belonged mainly to ST-15 and ST-5. Clade 4 comprised isolates from Egypt, Sudan and Palestine, with the majority belonging to ST-22 and ST-361. Clade 5 contained isolates from Egypt and belonged mainly to ST-97. The remaining isolates were in clade 6, of ST-239 and from Egypt and Morocco (n=1). This clade represents a *spa* t037/SCCmec III/MLST CC8 clone.

Figure 4. Phylogenetic tree based on core genome annotated by geographic origin, MLST CC, main *spa* and SCCmec types. 14 isolates, mostly from Egypt, lacked *mecA* and occurred predominantly in CC1 (n=5), CC15 (n=3) and CC30, CC8 and ST-80 (one isolate in each), ST-361 (n=2) or ST-5860 (n=1).

Table 1: MLST clonal complexes, *spa* types, and SCCmec types among the *S. aureus* isolates.

MLST CC	Strain Name	Geographical Origin	<i>spa</i> type	SCCmec type
CC1	3 (A), 3 (B), 23, 6 (B), 43, AA51, AA67, AA77	Egypt	t127	N/D
	R181, R180, AA1, AA78	UAE, Egypt	t127	N/D
	6 (A), AA69	Egypt	t127	N/D
	AA59, AA65, AA68	Egypt	t127	predicted as MSSA
	AA103, AA87	Egypt	t127	predicted as MSSA
CC15	15, 16	Egypt	t094	predicted as MSSA
	17	Egypt	unk	predicted as MSSA
CC22	41	Egypt	t13828	IV
	Gaza_MRSA_B62	Palestine	t223	IV
	AA18	Egypt	t223	IV
	AA5	Egypt	t3243	IV
	Gaza_MRSA_B04	Palestine	t790	IV
	40	Egypt	unk	IV
CC30	AA41	Egypt	t037	predicted as MSSA
	19	Egypt	t1504	N/D

CC5	AA30	Egypt	t304	IV
	14, AA76, AA80	Egypt	t688	N/D
	AA70	Egypt	t688	VI(4B)
CC8	12480433	Morocco	t008	IV
	LHI_Sa_30	Egypt	t008	N/D
	46	Egypt	t030	III
	50, AA101, AA13, AA14, AA22, AA23, AA27, AA31, AA33, AA46, AA52, AA55, AA57, AA60, AA61, AA62, AA63, AA64, AA91, AA92	Egypt	t037	III
	AA79	Egypt	t037	N/D
	AA93	Egypt	t037	predicted as MSSA
	AA29	Egypt	unk	III
CC97	AA36	Egypt	t267	N/D
	AA39, AA6	Egypt	t267	IV
	AA104	Egypt	t267	N/D
	AA8	Egypt	unk	IV
ST-80	2705432, 2705403, 2705405, 2705407, 2705409, 2705412	Tunisia,	t044	IV
		Kuwait,		
		Lebanon		
	AA45	Egypt	t044	IV
	2705431	Tunisia	t044	predicted as MSSA
	2705411	Lebanon	t131	IV
	AA2	Egypt	t416	IV
	AA3, AA4	Egypt	t416	N/D

420

421 **Table 2.** Virulence factors included in the *S. aureus* core genome.

VFclass	Virulence factors	Related genes
Adherence	Autolysin	<i>atl</i>
	Intercellular adhesin	<i>icaA</i>
		<i>icaD</i>
		<i>icaR</i>
Enzyme	Cysteine protease	<i>sspC</i>
	Thermonuclease	<i>nuc</i>
Immune evasion	Capsule	<i>cap5A</i>
		<i>cap8B</i>
		<i>cap5M</i>
		<i>cap8N</i>
		<i>capO</i>
Secretion system	Type VII secretion system	<i>esaB</i>
		<i>essA</i>
		<i>essB</i>
		<i>esxA</i>

422

423

References

1. Becker, K., Heilmann, C. & Peters, G. Coagulase-negative staphylococci. *Clinical Microbiology Reviews* **27**, 870–926 (2014).
2. Tong, S. Y. C., Davis, J. S., Eichenberger, E., Holland, T. L. & Fowler, V. G. Staphylococcus aureus infections: Epidemiology, pathophysiology, clinical manifestations, and management. *Clinical Microbiology Reviews* **28**, 603–661 (2015).
3. Stefani, S. *et al.* Meticillin-resistant Staphylococcus aureus (MRSA): Global epidemiology and harmonisation of typing methods. *International Journal of Antimicrobial Agents* vol. 39 273–282 (2012).
4. Chambers, H. F. & DeLeo, F. R. Waves of resistance: Staphylococcus aureus in the antibiotic era. *Nature Reviews Microbiology* vol. 7 629–641 (2009).
5. Froggatt, J. W., Johnston, J. L., Galetto, D. W. & Archer, G. L. Antimicrobial resistance in nosocomial isolates of Staphylococcus haemolyticus. *Antimicrobial Agents and Chemotherapy* **33**, 460–466 (1989).
6. Morrissey, I., Leakey, A. & Northwood, J. B. In vitro activity of ceftaroline and comparator antimicrobials against European and Middle East isolates from complicated skin and skin-structure infections collected in 2008-2009. *International Journal of Antimicrobial Agents* **40**, 227–234 (2012).
7. Maarouf, L., Omar, H., El-Nakeeb, M. & Abouelfetouh, A. Prevalence and mechanisms of linezolid resistance among staphylococcal clinical isolates from Egypt. *European Journal of Clinical Microbiology and Infectious Diseases* (2020) doi:10.1007/s10096-020-04045-w.

8. Steinig, E. J. *et al.* Evolution and global transmission of a multidrug-resistant, community-associated methicillin-resistant staphylococcus aureus lineage from the Indian subcontinent. *mBio* **10**, (2019).
9. Pobiega, M., Wójkowska-Mach, J. & Heczko, P. B. Typing of Staphylococcus aureus in order to determine the spread of drug resistant strains inside and outside hospital environment. *Przegląd epidemiologiczny* **67**, 435-438,539-542 (2013).
10. Lindsay, J. A. Evolution of Staphylococcus aureus and MRSA during outbreaks. *Infection, Genetics and Evolution* **21**, 548–553 (2014).
11. Funaki, T. *et al.* SCCmec typing of PVL-positive community-acquired Staphylococcus aureus (CA-MRSA) at a Japanese hospital. *Heliyon* **5**, e01415 (2019).
12. Sabri, I., Adwan, K., Essawi, T. A. & Farraj, M. A. Molecular characterization of methicillin-resistant Staphylococcus aureus isolates in three different Arab world countries . *European Journal of Microbiology and Immunology* **3**, 183–187 (2013).
13. Harris, S. R. *et al.* Whole-genome sequencing for analysis of an outbreak of meticillin-resistant Staphylococcus aureus: A descriptive study. *The Lancet Infectious Diseases* **13**, 130–136 (2013).
14. Eyre, D. W. *et al.* A pilot study of rapid benchtop sequencing of Staphylococcus aureus and Clostridium difficile for outbreak detection and surveillance. *BMJ Open* **2**, e001124 (2012).
15. Köser, C. U. *et al.* Rapid Whole-Genome Sequencing for Investigation of a Neonatal MRSA Outbreak. *New England Journal of Medicine* **366**, 2267–2275 (2012).
16. Tokajian, S. New epidemiology of Staphylococcus aureus infections in the Middle East. *Clinical Microbiology and Infection* vol. 20 624–628 (2014).

17. Borg, M. A. *et al.* Prevalence of methicillin-resistant *Staphylococcus aureus* (MRSA) in invasive isolates from southern and eastern Mediterranean countries on behalf of the ARMed Project members and collaborators. doi:10.1093/jac/dkm365.
18. Bhatta, D. R. *et al.* Association of Panton Valentine Leukocidin (PVL) genes with methicillin resistant *Staphylococcus aureus* (MRSA) in Western Nepal: A matter of concern for community infections (a hospital based prospective study). *BMC Infectious Diseases* **16**, (2016).
19. Falagas, M. E., Karageorgopoulos, D. E., Leptidis, J. & Korbila, I. P. MRSA in Africa: Filling the Global Map of Antimicrobial Resistance. *PLoS ONE* **8**, (2013).
20. Biber, A. *et al.* A typical hospital-acquired methicillin-resistant *staphylococcus aureus* clone is widespread in the community in the Gaza strip. *PLoS ONE* **7**, (2012).
21. Khalil, W., Hashwa, F., Shihabi, A. & Tokajian, S. Methicillin-resistant *Staphylococcus aureus* ST80-IV clone in children from Jordan. *Diagnostic Microbiology and Infectious Disease* **73**, 228–230 (2012).
22. Djoudi, F. *et al.* Panton-Valentine leukocidin positive sequence type 80 methicillin-resistant *Staphylococcus aureus* carrying a staphylococcal cassette chromosome *mec* type IVc is dominant in neonates and children in an Algiers hospital. *NEW MICROBIOLOGICA* vol. 36 (2013).
23. Enany, S., Yaoita, E., Yoshida, Y., Enany, M. & Yamamoto, T. Molecular characterization of Panton-Valentine leukocidin-positive community-acquired methicillin-resistant *Staphylococcus aureus* isolates in Egypt. *Microbiological Research* **165**, 152–162 (2010).
24. Abouelfetouh, A. The Status of Methicillin Resistance Among Egyptian *Staphylococcus aureus* Isolates: An Overview. *Infectious disorders drug targets* **17**, 67–69 (2017).

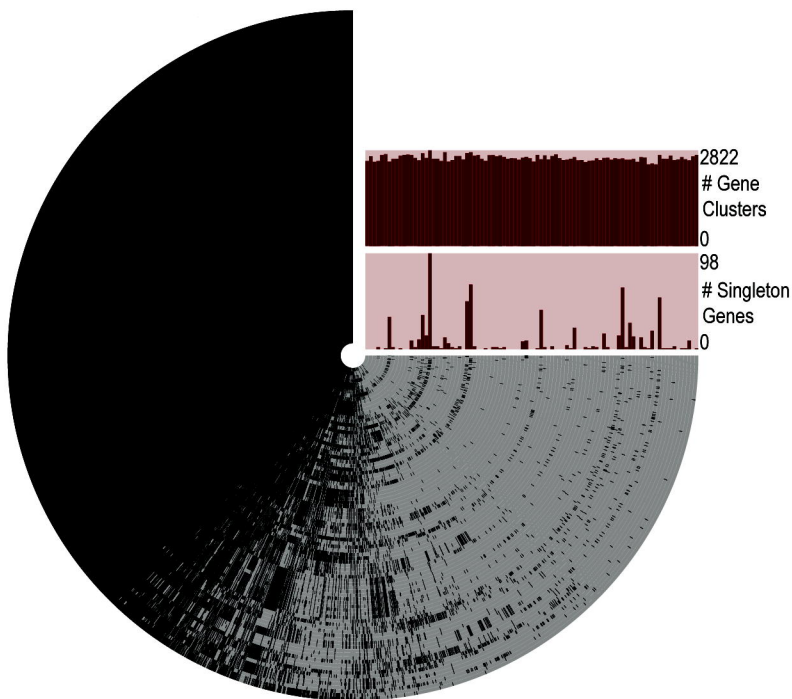
25. Deurenberg, R. H. *et al.* The molecular evolution of methicillin-resistant *Staphylococcus aureus*. *Clin Microbiol Infect* **13**, 222–235 (2007).
26. Guo, Y., Song, G., Sun, M., Wang, J. & Wang, Y. Prevalence and Therapies of Antibiotic-Resistance in *Staphylococcus aureus*. *Frontiers in Cellular and Infection Microbiology* **10**, (2020).
27. Köck, R. *et al.* Methicillin-resistant *Staphylococcus aureus* (MRSA): burden of disease and control challenges in Europe. *Euro Surveill* **15**, 19688 (2010).
28. Chongtrakool, P. *et al.* Staphylococcal cassette chromosome mec (SCCmec) typing of methicillin-resistant *Staphylococcus aureus* strains isolated in 11 Asian countries: a proposal for a new nomenclature for SCCmec elements. *Antimicrob Agents Chemother* **50**, 1001–1012 (2006).
29. Nickerson, E. K., West, T. E., Day, N. P. & Peacock, S. J. *Staphylococcus aureus* disease and drug resistance in resource-limited countries in south and east Asia. *Lancet Infect Dis* **9**, 130–135 (2009).
30. Kallen, A. J. *et al.* Health care-associated invasive MRSA infections, 2005-2008. *JAMA* **304**, 641–648 (2010).
31. Laupland, K. B. *et al.* The changing epidemiology of *Staphylococcus aureus* bloodstream infection: a multinational population-based surveillance study. *Clin Microbiol Infect* **19**, 465–471 (2013).
32. Yezli, S., Shibl, A. M., Livermore, D. M. & Memish, Z. A. Antimicrobial resistance among Gram-positive pathogens in Saudi Arabia. *J Chemother* **24**, 125–136 (2012).

33. John, J., George, S., Nori, S. R. C. & Nelson-Sathi, S. Phylogenomic Analysis Reveals the Evolutionary Route of Resistant Genes in *Staphylococcus aureus*. *Genome Biol Evol* **11**, 2917–2926 (2019).
34. Nguyen, M., Olson, R., Shukla, M., VanOeffelen, M. & Davis, J. J. Predicting antimicrobial resistance using conserved genes. *PLoS Comput Biol* **16**, e1008319 (2020).
35. Biswas, R. *et al.* Activity of the major staphylococcal autolysin Atl. *FEMS Microbiol Lett* **259**, 260–268 (2006).
36. O’Gara, J. P. *ica* and beyond: biofilm mechanisms and regulation in *Staphylococcus epidermidis* and *Staphylococcus aureus*. *FEMS Microbiol Lett* **270**, 179–188 (2007).
37. Jefferson, K. K., Pier, D. B., Goldmann, D. A. & Pier, G. B. The teicoplanin-associated locus regulator (TcaR) and the intercellular adhesin locus regulator (IcaR) are transcriptional inhibitors of the *ica* locus in *Staphylococcus aureus*. *J Bacteriol* **186**, 2449–2456 (2004).
38. Joo, H.-S. & Otto, M. Molecular basis of in vivo biofilm formation by bacterial pathogens. *Chem Biol* **19**, 1503–1513 (2012).
39. Khatoon, Z., McTiernan, C. D., Suuronen, E. J., Mah, T.-F. & Alarcon, E. I. Bacterial biofilm formation on implantable devices and approaches to its treatment and prevention. *Heliyon* **4**, e01067 (2018).
40. Pereira-Ribeiro, P. M. *et al.* Influence of antibiotics on biofilm formation by different clones of nosocomial *Staphylococcus haemolyticus*. *Future Microbiol* **14**, 789–799 (2019).
41. Donlan, R. M. Biofilms: microbial life on surfaces. *Emerg Infect Dis* **8**, 881–890 (2002).
42. Archer, N. K. *et al.* *Staphylococcus aureus* biofilms: properties, regulation, and roles in human disease. *Virulence* **2**, 445–459 (2011).

43. Morikawa, K. *et al.* Expression of a cryptic secondary sigma factor gene unveils natural competence for DNA transformation in *Staphylococcus aureus*. *PLoS Pathog* **8**, e1003003 (2012).
44. Lindsay, J. A. Genomic variation and evolution of *Staphylococcus aureus*. *Int J Med Microbiol* **300**, 98–103 (2010).
45. Langhanki, L. *et al.* In vivo competition and horizontal gene transfer among distinct *Staphylococcus aureus* lineages as major drivers for adaptational changes during long-term persistence in humans. *BMC Microbiol* **18**, 152 (2018).
46. Soliman, M. S. *et al.* Genomic Characterization of Methicillin-Resistant *Staphylococcus aureus* (MRSA) by High-Throughput Sequencing in a Tertiary Care Hospital. *Genes (Basel)* **11**, (2020).
47. Abou Shady, H. M., Bakr, A. E. A., Hashad, M. E. & Alzohairy, M. A. *Staphylococcus aureus* nasal carriage among outpatients attending primary health care centers: a comparative study of two cities in Saudi Arabia and Egypt. *Braz J Infect Dis* **19**, 68–76 (2015).
48. Alkharsah, K. R. *et al.* Comparative and molecular analysis of MRSA isolates from infection sites and carrier colonization sites. *Ann Clin Microbiol Antimicrob* **17**, 7 (2018).
49. Hadyeh, E., Azmi, K., Seir, R. A., Abdellatief, I. & Abdeen, Z. Molecular Characterization of Methicillin Resistant *Staphylococcus aureus* in West Bank-Palestine. *Front Public Health* **7**, 130 (2019).
50. Monecke, S. *et al.* Molecular Typing of ST239-MRSA-III From Diverse Geographic Locations and the Evolution of the SCCmec III Element During Its Intercontinental Spread. *Front Microbiol* **9**, 1436 (2018).

51. Alsequey, M. *et al.* Association between fluoroquinolone resistance and MRSA genotype in Alexandria, Egypt. *Scientific Reports* (2021).
52. Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).
53. Davis, J. J. *et al.* The PATRIC Bioinformatics Resource Center: expanding data and analysis capabilities. *Nucleic Acids Res* **48**, D606–D612 (2020).
54. Tatusova, T. *et al.* NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res.* **44**, 6614–6624 (2016).
55. Larsen, M. V. *et al.* Multilocus sequence typing of total-genome-sequenced bacteria. *J Clin Microbiol* **50**, 1355–1361 (2012).
56. Jolley, K. A., Bray, J. E. & Maiden, M. C. J. Open-access bacterial population genomics: BIGSdb software, the PubMLST.org website and their applications. *Wellcome Open Res* **3**, 124 (2018).
57. Bartels, M. D. *et al.* Comparing whole-genome sequencing with Sanger sequencing for spa typing of methicillin-resistant *Staphylococcus aureus*. *J Clin Microbiol* **52**, 4305–4308 (2014).
58. Kondo, Y. *et al.* Combination of multiplex PCRs for staphylococcal cassette chromosome mec type assignment: rapid identification system for mec, ccr, and major differences in junkyard regions. *Antimicrob Agents Chemother* **51**, 264–274 (2007).
59. International Working Group on the Classification of Staphylococcal Cassette Chromosome Elements (IWG-SCC). Classification of staphylococcal cassette chromosome mec (SCCmec): guidelines for reporting novel SCCmec elements. *Antimicrob Agents Chemother* **53**, 4961–4967 (2009).

60. Wattam, A. R. *et al.* Improvements to PATRIC, the all-bacterial Bioinformatics Database and Analysis Resource Center. *Nucleic Acids Res.* **45**, D535–D542 (2017).
61. Liu, B., Zheng, D., Jin, Q., Chen, L. & Yang, J. VFDB 2019: a comparative pathogenomic platform with an interactive web interface. *Nucleic Acids Res.* **47**, D687–D692 (2019).
62. Eren, A. M. *et al.* Anvi'o: an advanced analysis and visualization platform for 'omics data. *PeerJ* **3**, e1319 (2015).
63. Delmont, T. O. & Eren, A. M. Linking pangenomes and metagenomes: the Prochlorococcus metapangenome. *PeerJ* **6**, e4320 (2018).
64. Galperin, M. Y., Makarova, K. S., Wolf, Y. I. & Koonin, E. V. Expanded microbial genome coverage and improved protein family annotation in the COG database. *Nucleic Acids Res.* **43**, D261–269 (2015).
65. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–780 (2013).
66. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2 – Approximately Maximum-Likelihood Trees for Large Alignments. *PLoS ONE* **5**, e9490 (2010).
67. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments. *Nucleic Acids Research* **47**, W256–W259 (2019).

a**b**