



# Genome Sequencing of a Historic *Staphylococcus aureus* Collection Reveals New Enterotoxin Genes and Sheds Light on the Evolution and Genomic Organization of This Key Virulence Gene Family

Jo Dicks,<sup>a</sup> Jake D. Turnbull,<sup>a</sup> Julie Russell,<sup>a</sup> Julian Parkhill,<sup>b</sup> Sarah Alexander<sup>a</sup>

<sup>a</sup>Culture Collections, Public Health England, London, United Kingdom

<sup>b</sup>Department of Veterinary Medicine, University of Cambridge, Cambridge, United Kingdom

**ABSTRACT** We take advantage of a historic collection of 133 *Staphylococcus aureus* strains accessioned between 1924 and 2016, whose genomes have been long-read sequenced as part of a major National Collection of Type Cultures (NCTC) initiative, to conduct a gene family-wide computational analysis of enterotoxin genes. We identify two novel staphylococcal enterotoxin (pseudo)genes (*sel29p* and *sel30*), the former of which has not been observed in any contemporary strain to date. We provide further information on five additional enterotoxin genes or gene variants that either have recently entered the literature or for which the nomenclature or description is currently unclear (*selz*, *sel26*, *sel27*, *sel28*, and *ses-2p*). An examination of over 11,000 RefSeq genomes in search of wider support for these seven (pseudo)genes led to the identification of an additional three novel enterotoxin gene family members (*sel31*, *sel32*, and *sel33*) plus two new variants (*seh-2p* and *ses-3p*). We cast light on the genomic distribution of the enterotoxin genes, further defining their arrangement in gene clusters. Finally, we show that cooccurrence of enterotoxin genes is prevalent, with individual NCTC strains possessing as many as 18 enterotoxin genes and pseudogenes, and that clonal complex membership rather than time of isolation is the key factor in determining enterotoxin load.

**IMPORTANCE** *Staphylococcus aureus* strains pose a significant health risk to both human and animal populations. Key among this species' virulence factors is the staphylococcal enterotoxin gene family. Certain enterotoxin forms can induce a potentially life-threatening immune response, while others are implicated in less fatal though often severe conditions such as food poisoning. Genetic characterization of staphylococcal enterotoxin gene family members has steadily accumulated over recent decades, with over 20 genes now established in the literature. Despite the current wealth of knowledge on this important gene family, questions remain about the presence of additional enterotoxin genes and the genomic composition of family members. This study further expands knowledge of the staphylococcal enterotoxins while shedding light on their evolution over the last century.

**KEYWORDS** *Staphylococcus aureus*, enterotoxin gene family, genome analysis, National Collection of Type Cultures

*Staphylococcus aureus* is a Gram-positive, coccoid bacterium belonging to the *Firmicutes* phylum of mainly low-G+C bacteria. *S. aureus* is a common member of the human microbiota, with studies estimating approximately 20 to 30% of the population to be long-term carriers of *S. aureus* strains in the skin, nostrils, or female lower reproductive tract (1). In addition to its prevalence as a commensal organism of humans and animals,

**Citation** Dicks J, Turnbull JD, Russell J, Parkhill J, Alexander S. 2021. Genome sequencing of a historic *Staphylococcus aureus* collection reveals new enterotoxin genes and sheds light on the evolution and genomic organization of this key virulence gene family. *J Bacteriol* 203:e00587-20. <https://doi.org/10.1128/JB.00587-20>.

**Editor** Michael Y. Galperin, NCBI, NLM, National Institutes of Health

© Crown copyright 2021. The government of Australia, Canada, or the UK ("the Crown") owns the copyright interests of authors who are government employees. The [Crown Copyright](#) is not transferable.

Address correspondence to Sarah Alexander, [Sarah.Alexander@phe.gov.uk](mailto:Sarah.Alexander@phe.gov.uk).

**Received** 21 October 2020

**Accepted** 16 February 2021

**Accepted manuscript posted online** 1 March 2021

**Published** 21 April 2021

*S. aureus* is an important opportunistic pathogen. Strains can produce a variety of exotoxins, key among which are the staphylococcal enterotoxins (SEs), emetic toxins widely implicated in food poisoning. Gene family members are also associated with more severe, life-threatening conditions. For example, SEB is classified as a potential bioterrorism threat given its rapid and acute stimulation of the immune system, and it is also potentially implicated in the inducement of autoimmunity (2). Toxic shock syndrome (TSS) is a serious, and potentially fatal, condition with roughly half of cases denoted as menstruation associated and the remainder as non-menstruation associated. TSST-1, a protein very closely related to the SEs, gives rise to the majority of menstruation-associated TSS cases and approximately half of the nonmenstrual cases, with the remainder, ~25% in total, associated with SEB and SEC (3, 4).

The SE and TSST-1 proteins are superantigens (SAGs), immunomodulatory toxins that have the ability to stimulate large populations of T cells by interacting with the variable region of the  $\beta$ -chain (V $\beta$ ) of the T cell receptor. Structurally, SAGs are two-domain proteins characterized by a  $\beta$ -grasp domain and an OB-fold domain. The SE proteins are encoded by a family of genes related by their DNA sequence. The recent literature on the staphylococcal enterotoxin gene family encompasses 24 genes: *sea*, *seb*, *sec*, *sed*, *see*, *seg*, *seh*, *sei*, *selj*, *sek*, *sel*, *sem*, *sen*, *seo*, *sep*, *seq*, *ser*, *ses*, *set*, *selu*, *selv*, *selw* (formerly *selu2*), *selx*, and *sely*. The protein encoded by the closely related *tsst-1* gene was initially discovered independently by two groups as PEC (staphylococcal pyrogenic exotoxin C [5]) and SEF (staphylococcal enterotoxin F [6]) and later renamed TSST-1 upon agreement, hence the absence of the *sef* nomenclature within the SE gene list. Only genes whose proteins have demonstrated emetic activity are given the “*se*” prefix, with others designated “*sel*” for “staphylococcal enterotoxin-like.” The nomenclature used here is largely taken from the work of Fisher et al. (7).

In phylogenetic terms, the majority of the SE genes group closely with one another and with a number of *Streptococcus pyogenes* genes (8). The *selx* and *tsst-1* gene sequences group more distantly, with those of approximately 26 staphylococcal superantigen-like exoproteins (SSLs), which unlike the SAGs are immune invasion molecules and will not be considered here further. Despite its phylogenetic placement among the SSLs, *selx* is functionally similar to the SAG genes and is therefore referred to as an “SSL-like SAG” (8). Uniquely, SEIX is a single-domain SAG, lacking the OB-fold domain seen in all other staphylococcal SAGs to date.

The majority of the SAG genes are located on mobile genetic elements such as pathogenicity islands, prophages, and plasmids (7, 9). Clustering of the enterotoxin genes is also observed, most notably the *egc* cluster, which in a given strain can comprise up to seven genes and pseudogenes from a repertoire of nine (pseudo)gene forms (10). Consequently, while there is considerable variability with regard to the enterotoxin gene content between strains, the cooccurrence of the individual genes is highly nonrandom.

The National Collection of Type Cultures (NCTC) was founded in 1920 to address a recognized need for accumulating and disseminating information on human, animal, fungal, and plant pathogens. It is one of four culture collections operated by Public Health England as part of a globally recognized biological resource center, providing many thousands of historical and emerging strains to researchers and biomedical scientists worldwide. Recently, a Wellcome-funded initiative to sequence the genomes of ~3,000 NCTC strains was completed. Among the data sets developed were Pacific Bioscience (PacBio) long-read sequences and associated genome assemblies for 133 *Staphylococcus aureus* strains accessioned between 1924 and 2016, with at least one strain isolated prior to June 1924. NCTC strains are accessioned either proactively, based upon scientific requests from one of its nine past and current curators, or passively, deposited by members of the research community. While we think it unlikely that such a moderately sized data set would be representative—either geographically or temporally—of globally circulating *S. aureus* strains over the last century, the data set is nonetheless diverse, with 43 distinct sequence types (STs) each represented by one or more strains.

Here, we analyzed this new data set on a historic strain collection to answer questions on the number of staphylococcal enterotoxin genes and their genomic organization. We showed that each examined strain possessed between two and 18 SE genes. We identified seven putative SE genes outside our search list, four of which were not seen in NCTC strains accessioned after 1951, and one of which is the most prevalent enterotoxin-like sequence identified to date. We also examined the genomes of over 11,000 *Staphylococcus aureus* strains in the RefSeq database in order to gain support for this expanded SE gene repertoire. The RefSeq data set offered significant support for the newly identified genes while additionally presenting a further three SE genes and two gene variants. Collectively, the two data sets shed light on the genomic distribution of SE genes, further delineating six gene clusters and introducing a new one. Crucially, the NCTC data set enabled the examination of temporal patterns of enterotoxin birth and death to be made over a period of a century, showing a remarkable stability of gene content over this time, with all but one gene well represented in global sequence data sets. Finally, in accordance with this observed stability, analyzing the interrelationships between the NCTC strains showed that their clonal complex (CC) origins were more important than their time of isolation in determining their enterotoxin gene load.

## RESULTS

**The NCTC strains revealed novel enterotoxin-like sequences.** Using a profile hidden Markov model (pHMM) approach to hunt for DNA sequences within a set of genomes provides the opportunity to find novel SE genes that have yet to be formally characterized. We identified 825 SE- and 20 *tsst-1*-like sequences within the 133 *S. aureus* strains, with genomic coordinates and annotation details provided within Data Set S1 in the supplemental material. The 845 sequences clustered into 29 easily distinguishable gene-specific groups. In addition to finding 22 of the 25 expected SE/*tsst-1* gene-specific groups (reference sequences for each group are shown in Table S1; no copies of *see*, *ses*, or *set* were identified in any of the 133 strains), seven additional putative SE genes were identified, which we initially termed Gr1 to Gr7. Strikingly, 133 copies of the Gr1 sequence were found, one in every strain analyzed, spanning 43 distinct sequence types (including four not found in the *S. aureus* BIGSdb). While almost half (62 out of 133 copies) were likely to be pseudogenes due to premature stop codons or frameshift-inducing indels, seemingly intact versions of the coding sequence were seen across the accessioning period, with the most recent confirmed copy seen in strain NCTC 13434, isolated and accessioned in 2008.

All other groups (designated Gr2 to Gr7) consisted of between two and 16 members. The 16 copies of Gr2 were found in strains isolated between 1932 and 2008. Gr7 was less prevalent, with only five copies, but was observed over only a slightly reduced timespan, between 1938 and 1997. Despite its wide timespan, all five copies of Gr7 are likely to be pseudogenes. The other four genes were limited to the earlier strains, with Gr3/Gr4 (colocated), Gr5, and Gr6 most recently seen in strains accessioned in 1951, 1949, and 1948, respectively. Furthermore, both copies of Gr5 were presumed pseudogenes, likely the result of two small deletions. However, the prevalence of pseudogeny outside the cases of Gr1, Gr5, and Gr7 was generally less common, with fewer potential or likely pseudogenes in all other gene groups (see Data Set S1 for details).

**Origins of the putative enterotoxin-like sequences.** We searched for close sequence matches to each of the seven initially unidentified enterotoxin-like sequences in order to establish whether they had been observed previously. Overall, we found 882 high-scoring *megablast* hits to the GenBank nr/nt database (on 12 February 2020), with sequence-specific frequencies shown in Table S2a.

**(i) Gr1 is the near-ubiquitous chromosomal gene *seI26*.** The majority of the *megablast* hits, numbering 735 (83%, excluding 32 hits to NCTC genome sequences), were to the Gr1 nucleotide sequence at 93.24 to 100% sequence identity, including two full-length copies in a single strain (*S. aureus* strain ch22 chromosome [CP017807.1]) and all but 17 of which were to complete genomes or chromosomes. The 17 gene hits

were all annotated as enterotoxin-like W genes (e.g., 98.14% sequence identity to a purported *selw* gene identified in strain TD101 [KX655716.1]), though this gene was significantly different from the *egc* gene cluster *selw* gene used in our HMM search. It would appear that two distinct genes have been using the *selw* nomenclature: the *egc* gene formerly known as *selu2* and the seemingly ubiquitous (or at least highly prevalent) chromosomal gene previously seen in studies of human (11) and bovine (12) SEs. We will henceforth refer to this chromosomal gene as *sel26*, generally in line with the recommended nomenclature for SE genes (13) and similar to that used within reference 12, though currently lacking experimental confirmation to the best of our knowledge.

**(ii) Gr2 is the *orfX*-associated gene *selz*.** We found 40 hits to the Gr2 nucleotide sequence, all at 96.54 to 99.49% sequence identity. Fourteen hits were to gene sequences, the majority to a gene recently described as a staphylococcal enterotoxin-like Z gene (*selZ*) in strains of *Staphylococcus argenteus* (14). Interestingly, four hits (e.g., KT316803.1) were annotated as a staphylococcal cassette chromosome *mec* (SCC*mec*) element, a mobile genetic element implicated in broad-spectrum beta-lactam resistance via the *mecA* gene (15). An additional hit (U10927.2) appears to lie adjacent to an SCC*cap1* element, an SCC element with structural similarities to SCC*mec* but which instead harbors a type 1 capsular polysaccharide biosynthesis gene cluster (16). Further investigation of the annotation of U10927.2 shows that *selz* is the enterotoxin gene observed but unnamed in 2002 by Luong et al. (16). We searched for proximity of *selz* to *orfX* (which encodes an RlmH-type ribosomal methyltransferase) within each of the 16 strains with copies of Gr2/*selz*, as SCC elements are known to insert within the C terminus of this locus (17). We found that all copies of Gr2 were indeed in close proximity to *orfX* (between 2.5 kb and 38.6 kb), with two strains (NCTC 10399 and NCTC 10649) possessing an adjacent SCC*cap* element but none an SCC*mec* element. Six Gr2/*selz* copies were found in strains with sequence type (ST) 121, with others belonging to ST123, ST151, ST351, ST395, ST705, ST707, ST1254, and a novel sequence type. None of these strains belong to the six major identified clonal complex groups (CC1, CC5, CC8, CC22, CC30, and CC97) in this study.

**(iii) Gr3 and Gr4 are the clustered genes *sel27* and *sel28*.** The clustered Gr3 and Gr4 nucleotide sequences found 39 hits each, up to 98.41% and 99.12% sequence identity, respectively. Of particular note were strong hits for each gene to a gene cluster pathogenicity island in strain 364P and to two recently identified enterotoxin genes annotated as *Sel27* and *Sel28* in strains SJTU F20365 (MF370878.1 [18]), 86, 72, 50, SG19, SG16, SG13, SG11, SG09, SG05-2, SG05-1, SG04, and SG01. The oldest confirmed NCTC strain identified as carrying these two genes, NCTC 5664, was isolated in 1936 and the putative youngest, NCTC 8765, during or prior to 1951. Five of the seven NCTC strains possessing the two genes belong to ST9 (CC1) and the remainder to ST350.

**(iv) Gr5 is the (pseudo)gene *sel29p*.** We failed to find any highly similar hits to the Gr5 nucleotide sequence. The two strains possessing this gene, NCTC 6966 and NCTC 7856, were accessioned in 1945 and 1949, respectively; both belong to ST890; and the two sequences appear to be pseudogenes. We propose that this gene be referred to henceforth as *sel29p*.

**(v) Gr6 is the plasmid gene *sel30*.** The Gr6 nucleotide sequence produced just eight high-scoring hits, all to complete plasmid sequences and differing at most by one nucleotide substitution. The absence of any hits to annotated coding sequences meant it was difficult to form any further conclusions about the origins of this gene, other than its clear plasmid location. The two NCTC strains possessing the Gr6 gene were both accessioned in the 1940s and belonged to ST5 and ST1021. We propose that this gene be referred to henceforth as *sel30*.

**(vi) Gr7 is the *orfX*-associated pseudogene *ses-2p*, a variant of *ses* clustered with *seh*.** Finally, the Gr7 nucleotide sequence, present in all five NCTC copies in close proximity to an *seh* gene sequence, found 21 hits ranging between 93.51% and 100% sequence identity. While the majority were to complete genome or chromosome sequences, two hits (EU272079.1 and KX690110.1) were—similarly to Gr2—to insertion sites of SCC*mec* elements. Further investigation of these hits identified previous reports

of a partial enterotoxin gene with sequence similarity to *seo*, in close proximity to an *seh* gene and associated with SCCmec type IV element insertion (19, 20). Similar to our analysis of Gr2, we searched for *orfX* genes within the five strains with copies of both Gr7 and *seh*, plus the single strain (NCTC 13435) possessing a presumed pseudogenized version of *seh* but no copy of Gr7. We found that all five copies of Gr7 and six copies of *seh/sehp* were close to an *orfX* gene (between 17.5 kb and 43.1 kb). However, only one copy of Gr7 (in NCTC 13297) was adjacent to an SCC element (in this case likely an SCCfus element) while the *sehp* copy in NCTC 13435 was adjacent to an SCCmec type IV element. Further sequence analysis showed that while the 3' region of the Gr7 gene (~300 bp) was highly similar to the corresponding region of *ses* (rather than *seo*), the 5' region of approximately 30 bp was ~40 bp shorter and dissimilar at the sequence level, with the intervening region showing a moderate level of sequence similarity interrupted by several presumed mutations. All copies of Gr7 in NCTC strains therefore look to be (only partially) truncated pseudogenes, containing several premature stop codons, though it is uncertain whether Gr7 was ever a functional gene. Four of the NCTC strains with the Gr7 sequence are ST10, with the remaining strain ST1. We propose that this pseudogene be referred to henceforth as *ses-2p*.

**Most novel enterotoxin-like gene and pseudogene sequences are also observed within RefSeq genomes.** We investigated the numbers of high-scoring hits for each of the seven novel or recently identified (pseudo)genes within 11,351 *Staphylococcus aureus* genome sequences within the RefSeq database (on 11 May 2020). Of the 64,281 total hits to staphylococcal enterotoxin-like sequences, 12,505 were to six of these seven genes. Only Gr5/*sel29p* failed to find hits in strains other than the two NCTC strains (which were present within the RefSeq database). Table S2a in the supplemental material shows that the relative frequencies of the other six genes within the NCTC data set are mirrored within the nr/nt and RefSeq databases. This suggests to a certain extent that the frequencies of SE genes within the NCTC data set are indicative of their frequencies within larger data sets, notwithstanding the absence of three SE genes from our data set. Although we did find instances of *see*, *ses*, and *set* within the sizeable RefSeq data set (2, 39, and 36 copies, respectively), their low frequencies suggest these may be relatively rare genes, or certainly within strains whose genomes have been sequenced thus far.

**The RefSeq database harbors a further cache of novel enterotoxin-like sequences.** We classified 11,026 of the 11,351 RefSeq genomes into 468 distinct sequence types (see Data Set S2), 39 of them putative new STs. Analysis of all genomes led to the identification of a further three putative staphylococcal enterotoxin genes plus two additional gene variants. In total, 258 sequences grouped into five distinct sets which we initially termed Gr8 to Gr12. Similar to the NCTC-derived sequences, we searched the GenBank nr/nt database using BLAST to determine further information on the origins of these groups (see Table S2b).

**(i) Gr8 and Gr9 are the clustered genes *sel31* and *sel32*.** Fifteen copies of the Gr8 sequence in the RefSeq genomes were to a likely functional (based on its amino acid translation) gene which we term *sel31*. Of the 18 Gr9 sequences, 15 directly neighbored *sel31*. We term this new gene *sel32* and hence delineate a new enterotoxin gene cluster. Fifteen of the 18 RefSeq genomes could be classified as four sequence types (ST1, ST121, ST97, and ST508) from three clonal complexes (CC1, CC45, and CC97); see Data Set S3 for details of these strains. Each gene found three identical BLAST hits to the GenBank nr/nt database, from the same genomic sources, two of which were to plasmid genomes, highlighting the likely origin of the gene pair.

**(ii) Gr10 is the egc cluster gene *sel33*, a recombinant of *selw* and *sen*.** A single copy of a new recombinant derivative of egc gene cluster genes *selw* and *sen* was identified in strain BSAC1477 from the BSAC Resistance Surveillance Project (GenBank accession no. [NZ\\_FGM101000018.1](https://www.ncbi.nlm.nih.gov/nuccore/NZ_FGM101000018.1); <http://bsacsurv.org/>), which we term *sel33*. This strain was isolated in or after 2001 and derives from CC22. Of the five sequences investigated here, only *sel33* failed to find any similar sequences within the nr/nt database, suggesting this recombination to be a rare occurrence.



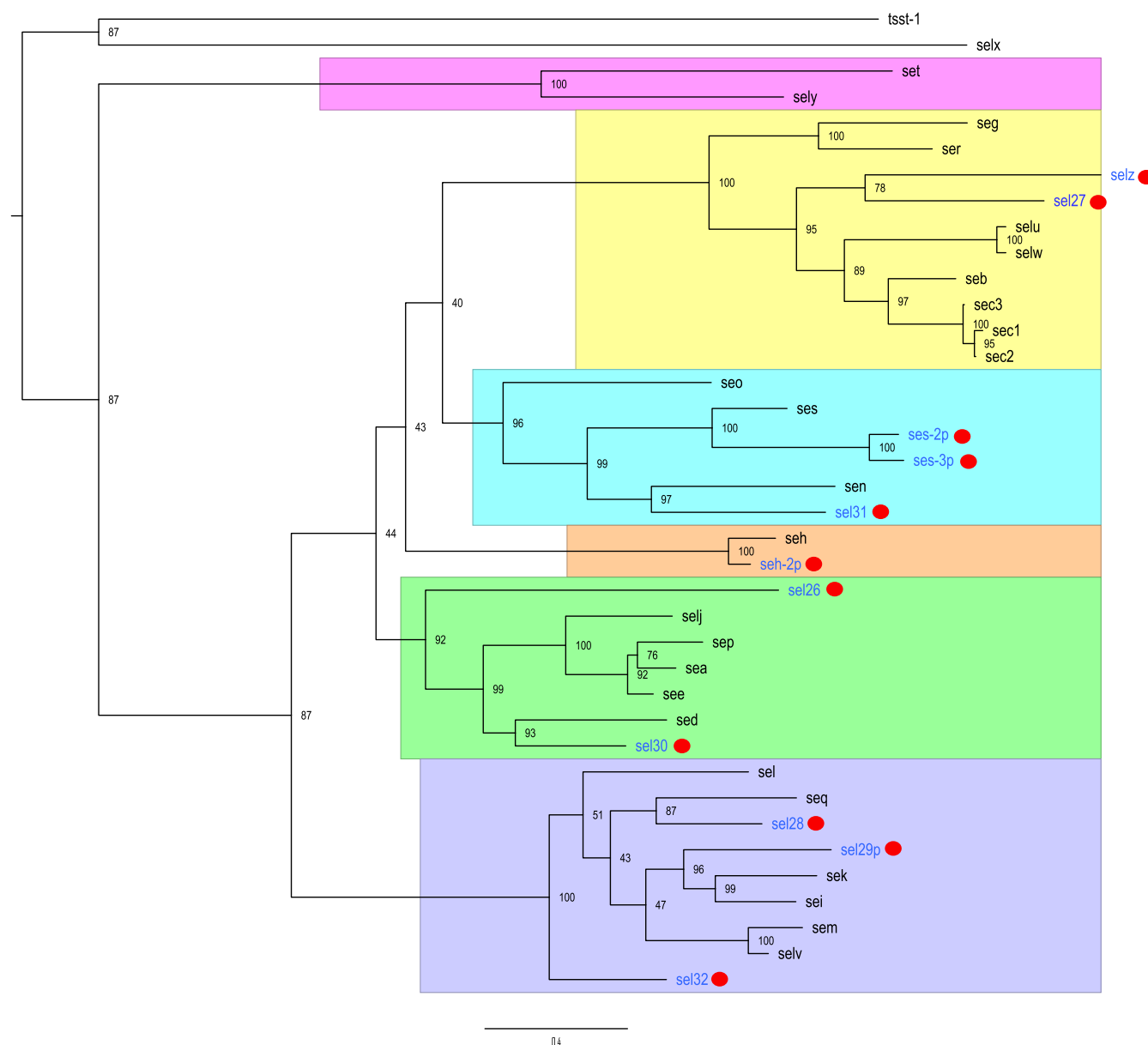
**(iii) Gr11 is the variant of *seh*, *seh-2p*, strongly associated with SCCmec type IV elements.** The RefSeq data set contained 30 copies of Gr11, a pseudogenized or truncated form of *seh* which we refer to henceforth as *seh-2p*. The sequence of *seh-2p* likely possesses two single nucleotide deletions relative to the canonical form of *seh*. All but two classified genomes derive from ST80, with one each of ST4563 and a novel sequence type, both differing in only a single multilocus sequence typing (MLST) allele from ST80 (both within the *glpF* gene). Notably, one of the ST80 strains is NCTC 13435, which, unlike the five copies of *seh* in the NCTC strains, lacked a neighboring *ses-2p* sequence. Similarly, none of the 30 copies of *seh-2p* within the RefSeq data set has a neighboring enterotoxin-like sequence, and 29 show significant evidence of a neighboring SCCmecVc(2B) element, with 27 genomes showing all SCCmec genes spread over one to four contigs and two additional genomes showing partial SCCmec matches (both including *mecA* presence). In contrast, 219 of the 224 copies of *seh* possessed an adjacent *ses-2p* or *ses-3p* (see below) sequence, with 29 and 190 copies, respectively.

**(iv) Gr12 is the pseudogene *ses-3p*, a further variant of *ses* linked to *seh*.** We saw above that five copies of *ses-2p* were found within the NCTC genomes adjacent to copies of *seh*, and 29 such gene pairs were also observed within the RefSeq genomes. The Gr12 group of 194 sequences was found to constitute a second, distinct variant of *ses* which we term *ses-3p* given its likely pseudogene status. This new variant is highly similar to *ses-2p* except for a divergent 5' end. As noted above, most copies of *ses-3p* were found adjacent to *seh*, with only four of 194 instances lacking a neighboring enterotoxin sequence. All but four classified genomes were found to be members of CC1 (ST1, ST81, ST474, ST1207, ST2764, ST3248, ST3497, and a novel ST), with the remainder from ST182 and ST944, sequence types highly distinct from CC1 but differing from one another in a single allele.

**The novel enterotoxin gene sequences are spread across the SE phylogeny.** A phylogenetic tree (Fig. 1) was estimated from the amino acid sequences of 11 of the 12 putative novel or recently identified SE genes (or "repaired" amino acid sequences in the cases of *sel29p*, *ses-2p*, and *ses-3p*, and a "short" sequence truncated by a premature stop codon in the case of *seh-2p*), alongside sequences of the established SE genes that were used in the pHMM search process. The tree shows the 11 sequences to group across the SE gene tree: *sel28*, *sel29p*, and *sel32* with *sei*, *sek*, *sel*, *sem*, *seq*, and *selv*; *sel26* and *sel30* with *sea*, *sed*, *see*, *selj*, and *sep*; *seh-2p* with *seh*; *ses-2p*, *ses-3p*, and *sel31* with *sen*, *seo*, and *ses*; *selz* and *sel27* with *seb*, *sec*, *seg*, *ser*, *selu*, and *selw*. The tree groupings of the established SE genes remained largely consistent with earlier analyses (e.g., reference 8) following the addition of the new gene family members. Note, however, that the amino acid sequence of *sel33* was omitted from the tree. As it is a recombinant gene derived from genes in two distinct clades (*selw* in the yellow clade and *sen* in the cyan clade in Fig. 1), its inclusion distorts the topology of the resulting tree. This contrasts with *selv*, which derives from two genes within the same clade (*sem* and *sei*, purple clade in Fig. 1).

**There are at least seven staphylococcal enterotoxin gene clusters.** SE genes are known to sometimes collocate with others, often on plasmids or pathogenicity islands. The most striking example of this phenomenon is the *egc* gene cluster. The full characterization of this operon has taken place in a stepwise fashion. The initial discovery of the *seg* and *sei* genes (21) was followed by identification of *sem*, *sen*, and *seo* in the neighboring genomic regions, along with evidence of their cotranscription, while two pseudogenes (*φent1* and *φent2*) were found between *sei* and *sen* (22). A sixth gene, *selu*, thought to be the product of deletions within *φent1* and *φent2*, was identified later (23). Most recently, *selw* (formerly *selu2*) and *selv* were identified (10). While the nucleotide sequence of *selw* was highly similar to that of *selu*, the main difference being a 15-bp deletion in the former compared to the latter, and thought to be the result of a different mutation of the *φent1* and *φent2* sequences from that hypothesized in *selu*, *selv* was found to be the product of a recombination of *sem* and *sei*.

The *egc* gene cluster appears to be highly prevalent within *S. aureus* genomes (14, 24). In this study, 59 *egc* gene clusters were found in 58 of the 133 *S. aureus* strains



**FIG 1** Phylogeny of the staphylococcal enterotoxin genes. Maximum likelihood (ML) phylogenetic tree of 11 of the 12 new staphylococcal enterotoxin gene family members (gene names shown in blue text and with adjacent red circle; *sel33* is not shown as its between-clade recombinant origin distorts the tree topology) identified in the NCTC and RefSeq strain sets, along with reference sequences for 24 previously identified SE genes (including three variants of *sec*) plus *tsst-1*. Compact gene groups (clades) are highlighted as colored blocks. The tree was estimated with IQ-TREE (27) using the VT+F+R4 amino acid substitution model, maximum log-likelihood =  $-14,507.9726$ , and with the *tsst-1/selx* clade used as an outgroup. One thousand ultrafast bootstraps were performed, with percentages of bootstrapped trees supporting the ML tree shown at each internal node. The tree was further annotated by clade with FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>).

(43.6%), all but one seemingly complete. One of these 58 strains (NCTC 7972) possessed two gene clusters, with one of them the only case observed in this strain set of a large gap between any two *egc* genes, 13.6 kb between *seo* and *sem*. In NCTC 11963, the only clear case of an incomplete *egc* gene cluster, we identified genes *seg* and *seo* on two separate genomic contigs, such that the omission of intervening genes may be the result of an incomplete genome assembly rather than their absence from the genome. Further investigation of the raw sequencing reads overlapping this region showed the central region of the *egc* gene cluster to suffer from a very low read coverage but nonetheless to offer support for the existence of a full gene cluster most likely of type OMIUNG (read 27628 runs from the middle of *seo* to the middle of *seg* and is

**TABLE 1** *egc* gene cluster forms within 133 NCTC *Staphylococcus aureus* genomes

Cluster composition	No. of copies
OMIUNG	39
OMIWNG	15
OMIUN	3
OVUNG	1
Unconfirmed (O and G only)	1
OVWNG	0
OMI33G	0
Total	59

highly similar to the corresponding sequence of NCTC 2669, albeit with a single large run of T's breaking the alignment, presumably a sequencing artifact).

The NCTC data set suggests the  $\phi$ ent1 and  $\phi$ ent2 pseudogenes should no longer be considered entities distinct from *selu* and *selw*. While all four pseudogenized copies of *selw* and one of four pseudogenized copies of *selu* possess a single nucleotide frameshift (a run of 6 A's increased to 7 A's at positions 365 in *selw* and 380 in *selu*) that would lead to a two-open reading frame (ORF) prediction similar to  $\phi$ ent1 and  $\phi$ ent2, the underlying nucleotide sequences are clearly merely a minor change to *selu* or *selw*. Furthermore, none of these sequences possesses the 69-bp deletion (relative to *selw*) observed in strain A900322, from which  $\phi$ ent1 and  $\phi$ ent2 were first defined (22), nor could we infer this deletion from any other strain sequence within the GenBank database. We feel that it is therefore more appropriate going forward to refer to *selu* or *selw* and their pseudogenes only. Given that the earliest copies of these genes within the NCTC collection (NCTC 2669 from 1928 for *selu* and NCTC 6134 from 1941 for *selw*) appear to be full-length, functional copies, the historical data would also support this view.

Frequencies of the distinct *egc* gene cluster arrangements identified in the NCTC strains are given in Table 1, showing that OMIUNG (i.e., the gene order *seo-sem-sei-selu-sen-seg*) and its close variant OMIWNG (including "minor" pseudogenes of all six genes) are the predominant gene cluster variants, seen in this strain set in a ratio of 2.6:1. As well as three strains isolated and accessioned in the 1930s or 1940s possessing an OMIUN variant (i.e., apparent absence of the *seg* gene), we see a recent strain (NCTC 13373, accessioned in 2005 and equivalent to ATCC 43300, a clinical isolate from Kansas) with the OVUNG form, the potentially rare gene *selv* the result of a recombination between *sem* and *sei*. A comparison of the sequence of *selv* in NCTC 13373 to the canonical form in strain A900624 (10) from the French National Reference Center for Staphylococci (see Fig. S1 in the supplemental material) shows evidence for the recombination between *sem* and *sei* in NCTC 13373 having occurred slightly closer to the 5' end of the sequence, though both events clearly took place within a central sequence region highly similar between the two progenitor genes. That observation, together with a high number of single nucleotide differences between the NCTC 13373 and *selv* reference sequences plus the alternative OVWNG form of the A900624 gene cluster, indicates that the two genes likely arose from two distinct recombination events. Note that the OVWNG form and the OMI33G form (i.e., containing the novel *sel33* recombinant of *selw* and *sen*), which we discovered in strain BSAC1477, were not observed within the NCTC data set.

Five additional gene clusters were observed within the NCTC strains, as shown in Table 2. All but one instance of the 122 gene clusters, the broken *egc* cluster mentioned above, were found in intact gene cluster form. Interestingly, the *sel27-sel28* gene cluster, found within 7 strains accessioned in the 1930s to 1950s, was seen to be located close to the *egc* gene cluster (2 of the 3 OMIUN strains and 5 of the 14 OMIWNG strains). The distance of the *egc* cluster to the *sel27-sel28* gene cluster ranged between 17,465 and 39,464 bp, with an average of approximately 27.6 kb.



**TABLE 2** *Staphylococcus* enterotoxin gene clusters within 133 NCTC *Staphylococcus aureus* genomes

Cluster composition	No. of copies
egc	59
sek-seq	23
sec-sel	19
sed-selj-ser	9
sel27-sel28	7
seh-ses-2p	5
sel31-sel32	0
Total	122

While the *sek-seq* and *seh-ses-2p* gene clusters are comprised of genes within the same clade in Fig. 1 (the cyan and orange groups are often referred to as a single clade elsewhere), the remaining four clusters contain genes spanning two or even three clades. In particular, the two most prevalent *egc* gene cluster arrangements (OMIUNG and OMIWNG), which account for 54 of its 59 copies, possess two genes from three main SE clades (shaded yellow, cyan, and purple in Fig. 1). It has been speculated that this divergence of the *egc* cluster genes may indicate the cluster's role as the progenitor of the majority of SE genes in *Staphylococcus aureus* (22). The seventh gene cluster, *sel31-sel32*, identified in 15 RefSeq genome sequences from three clonal complexes, was not observed within the NCTC strains.

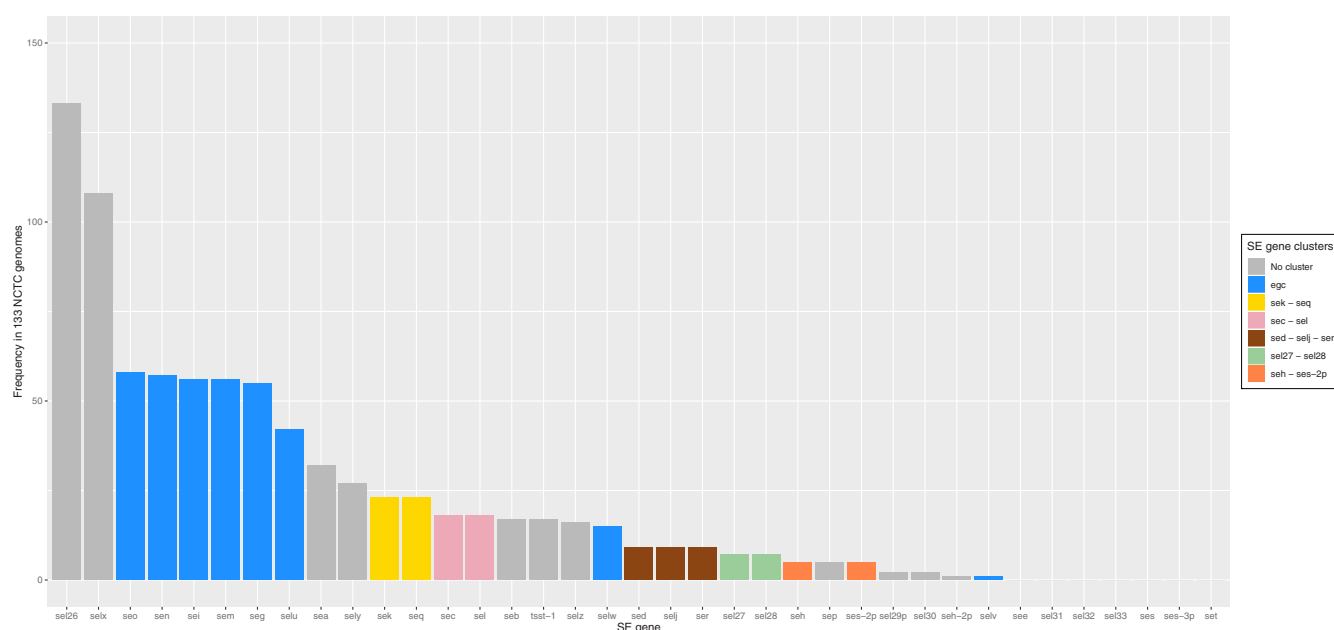
***Staphylococcus aureus* strains can possess many SE genes.** The frequencies of the 37 *SE/tsst-1* gene groups, including the 12 additional enterotoxin-like sequences, within the 133 NCTC strains are shown in Fig. 2. Data Set S1 further shows the numbers of each gene identified within each strain. We see that the chromosomal genes *sel26* and *selx* are most common, with lesser frequencies of genes present on mobile genetic elements. Further work would be necessary to determine whether these frequencies were representative of the population as a whole or whether they were biased by sampling and temporal effects.

We found 10 cases of a strain possessing two copies of the same gene and a single case of three gene copies (see Data Set S1). Notable examples include NCTC 7415, in which we found both three copies of *tsst-1* and two copies of the *sec-sel* gene cluster, and NCTC 7972, which harbored two intact *egc* gene clusters. In the former case, the likelihood of two of the three *tsst-1* copies and one of the *sec* copies being pseudogenes within a 20-kbp region of a single contig may contribute to this finding.

Consistent with other studies such as that of Varshney et al. (24), individual strains were found to possess numerous *SE/tsst-1* genes, with a range of 2 to 18 genes per strain and a mean of 6.35 (median of 6). While this value is slightly higher than the average 5 SE genes per strain seen in reference 24, that prior study had looked at fewer genes, 19 of the 37 genes examined here, which may have led to the lower gene counts.

**Associations between unclustered SE genes.** Unlike the 19 SE genes involved in gene clusters, and which we discussed above, the *sea*, *seb*, *sep*, *selx*, *sely*, and *tsst-1* genes are not clustered within this data set in a conventional form, as are neither the newly identified *selz*, *sel26*, and *sel30* genes nor the *sel29p* and *seh-2p* pseudogenes. Nevertheless, both positive and negative associations between these and other SE genes may still exist, likely the result of enterotoxin gene copresence on plasmids, prophages, pathogenicity islands, and other mobile genomic islands (9). Figure S2a shows a heat map of Pearson correlation coefficients of gene presence/absence for all gene pairs (with the exception of *sel26*, which is always present), with genes arranged so that they are close to other genes with which they show the greatest associations.

The six gene clusters identified in the section above are easily apparent as either single or sets of large red circles. Additional positive and negative associations are also apparent. Notable positive associations include *tsst-1* with the *sec-sel* gene cluster, *seb*

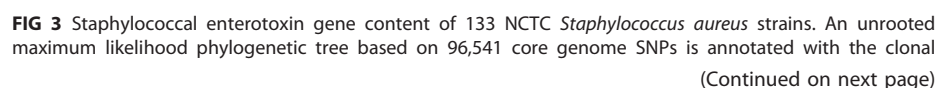


**FIG 2** Staphylococcal enterotoxin gene presence. Frequencies of 36 staphylococcal enterotoxin genes (or putative genes/pseudogenes) plus *tsst-1* in 133 NCTC *S. aureus* strains. Membership of one of the six gene clusters present in this data set is indicated by a color code.

with the *sek-seq* gene cluster, *sely* with *selz*, and *selw* with the *sel27-sel28* gene cluster. The former two associations are previously noted and likely the products of copresence on SaPI<sub>m1/n1</sub> (or SaPI<sub>bov1</sub>) and SaPI<sub>3</sub> pathogenicity islands, respectively. Examination of the relevant gene coordinates shows that the associations are without exception underpinned by copresence on the same contig, though not by colocation. The most compact examples are the 11 cases of *seb/sek-seq*, which are approximately 11 kb apart in all strains. However, to the best of our knowledge the latter two associations have not been observed prior to this study. *sely* has previously been seen only on the chromosome, so its association with the *orfX*-associated *selz* gene could be down to chance alone. The *selw/sel27-sel28* association is likely to involve particular forms of the  $\nu$ Sa $\beta$  genomic island, given the known presence of the *egc* cluster on this mobile element (25).

Notable negative associations include *selx* with the OMIUNG form of the *egc* gene cluster and the prophage-carried *sea/sek-seq* gene cluster combination with OMIUNG and *sely*. Although Varshney et al. (24) note an absence of *egc* gene clusters within *seb*<sup>+</sup> strains, we observe a more complex pattern in this data set. Examining the 17 *seb*<sup>+</sup> strains identified here, we find that 7 of the 8 strains isolated in or prior to 1949 carried the *egc* gene cluster whereas none of the 9 strains believed to derive from the 1950s onward were found to possess it. Furthermore, no strains with *sel27* and *sel28* harbored a *seb* gene. Taken together, these observations suggest that different SE gene combinations have circulated within the *S. aureus* population, some of which were restricted to particular time frames, or perhaps to different parts of the *S. aureus* population. Figure S2b shows a depiction of the associations observed in this study and/or described in the work of Argudin et al. (9).

**Phylogenetic and temporal patterns of staphylococcal enterotoxin genes.** We estimated a phylogenetic tree of the 133 NCTC *S. aureus* strains using HarvestTools (26) and IQ-TREE (27), based on 96,541 core genome single nucleotide polymorphisms (SNPs), and annotated it with SE/*tsst-1* gene content and clonal complex group. It is immediately clear from Fig. 3 that certain SE gene combinations are restricted to particular clades within the tree. For example, only CC1, CC5, and CC22 strains harbor the OMIUNG form of the *egc* gene cluster in this data set. Indeed, the observed patterns of gene content explain many of the associations between genes described in the



sections above. For example, we see that *selx* is completely absent from CC30 (purple strip in Fig. 3), a monophyletic group in which the OMIUNG form of the *egc* gene cluster is highly prevalent, thereby explaining the strong negative association between *selx* and *selu*.

Most clades cover a broad timespan (e.g., confirmed isolation periods of at least 1941 to 1997 for CC1, 1948 to 1988 for CC5, 1932 to 2003 for CC8, and 1928 to 2003 for CC30) with a potential span of 1933 to 1985 for CC97 (though presently confirmed only up to 1954) and only CC22 represented here by a narrow group of strains (1990 to 2005). While Fig. 3 suggests gene content to be highly correlated with clonal complex, the mean number of enterotoxin genes per strain was found to be higher in strains isolated within the 1920s to 1940s (7.60) than in the 1950s to 2010s (5.65), irrespective of clonal complex membership. We examined the relationship between the number of enterotoxin genes per strain with clonal complex and year of isolation by fitting a generalized linear model to the data, collapsing gene clusters to single observations as described in Materials and Methods. We found that for the 90 NCTC strains with clonal complex designations (see Fig. S4 in the supplemental material for plots of the data), the number of enterotoxin genes/gene clusters was strongly associated with clonal complex identity ( $P < 0.05$  for three of the five factors compared to CC1) but neither with year of isolation ( $P = 0.274$ ) nor with the interaction between clonal complex and time ( $P > 0.425$  for all factor interactions). This suggests that enterotoxin gene content within clonal complexes has remained stable across the century of strain isolation and that the putative temporal difference in gene content described above may be due to sampling effects, with a higher frequency of strains harboring the *egc* gene cluster isolated during the earlier period (60% versus 34%).

HarvestTools Gingr plots of the core genome SNP alleles alongside the phylogenetic tree (see Fig. S3 for an example, with NCTC 1803 as a reference genome) also indicate that horizontal transmission between disparate *S. aureus* clades has taken place. Consequently, staphylococcal enterotoxin gene content may also be influenced by horizontal processes in addition to clonal expansion. In future, it would be interesting to analyze whether, for example, the two cases of the *sel27-sel28* gene cluster outside CC1 (uppermost green circles in Fig. 3) were due to horizontal transfer of the pathogenicity island on which they are located.

## DISCUSSION

We analyzed a strain set of 133 *Staphylococcus aureus* strains from the UK National Collection of Type Cultures, with a particular goal of understanding the complement of enterotoxin genes captured within, thereby further enhancing the utility of the strains for the benefit of the research community. While we did not initially anticipate uncovering any potential novel genes, particularly given the size of the data set and the lack of an enterotoxin-focused strategy for its collection, the use of a pHMM profile approach allowed us to identify new sequences that we hope will be investigated further by researchers in this area. The relative ease with which we found putative enterotoxins first within the NCTC data set, and subsequently within the RefSeq database, leads us to speculate as to whether there might yet be other enterotoxin genes left for others to uncover. While the NCTC and RefSeq data sets encompass a sizeable proportion of the global diversity of *S. aureus* strains, with 43 and 468 distinct sequence types represented, respectively, they will not have captured the full range. Consequently, as-yet-unidentified enterotoxin genes may still be present in strains whose genomes are currently outside the reach of strain and sequence collections.

### FIG 3 Legend (Continued)

complex (color strip adjacent to strain names: CC1, blue; CC5, gold; CC8, red; CC22, green; CC30, purple; CC97, orange) and *SE/tst-1* gene content (established genes as squares and recent/novel genes as circles). Gene presence is colored according to the scheme in Fig. 2, so that membership of a common gene cluster can be identified easily. SNPs were called using HarvestTools. The tree was estimated with IQ-TREE using the SYM+ASC+R3 nucleotide substitution model, with a maximum log-likelihood = -1,022,309.6472. The figure was generated using the iTOL web server (49).

Our study has also added to the understanding of the genomic organization of the enterotoxin genes, particularly via gene clusters carried by mobile genetic elements. Gene families harbored by bacterial genomes are presented with an array of strategies that enable them to thrive and mobilize. The staphylococcal enterotoxin gene family has indeed shown it is capable of exploiting many of these routes, from use of plasmids, prophages, pathogenicity islands, and genomic islands in addition to stable chromosomal inheritance. Despite the consequent stability of the SE genes over the past century, in general large gene families appear to be rare in many bacteria. A study of the sequenced genomes of species including *Escherichia coli*, *Streptococcus pyogenes*, and *Chlamydomphila pneumoniae* showed limited numbers of gene families of size 20 or over, with only ~10 such gene families in *S. aureus* strains Mu50, MW2, and N315 (28). A more recent study found greater variation in the number of gene families between strains, but again the number of duplicated genes was relatively limited, with a maximum of 190 duplicates for 84 gene families across 473 strains (29). Furthermore, the majority of duplicated genes in the latter study were thought to have a phage origin. The findings here are consistent with this observation.

The *egc* gene cluster has clearly been a key component in the expansion of the SE gene family. Interestingly, the genomic island harboring the *egc* gene cluster was recently found, similarly to pathogenicity islands, to be capable of mobilization due to a temperate bacteriophage (30). The *egc* is the only SE gene cluster to date that has been shown to have produced novel recombinant SE genes, and from the growing number of components on offer, six distinct combinations were observed in the strains analyzed here. As mentioned earlier, the *egc* gene cluster has been mooted as a putative SE nursery, whereby the observed genetic diversity has been generated by the processes of tandem duplication and subsequent divergence (22). The level of variation observed in this study would seem consistent with that view. Looking at the seven gene clusters in Table 2, all but one (*sek-seq*) has members from two or more of the shaded clades in Fig. 1. In future, it might be illuminating to carry out a dating analysis of SE gene sequences to see if the results can help us to understand how these structures might have evolved. For example, the constituents of the *sec-sel* and *se27-se28* clusters derive from the same two clades, and further, these two clades are two of the three clades from which the *egc* genes all derive. It would be interesting to determine whether these common features are due to a (partially) shared inheritance or whether they are coincident with independent origins.

The cases of *seh* (*ses-2p*) and *selz* genes also indicate that transposition of genes to insertion sites otherwise used by elements such as staphylococcal cassette chromosomes may be an additional strategy for gene survival and proliferation. Interestingly, Luong et al. (16) and Noto et al. (19) have suggested that enterotoxin insertion at this genomic location may have been implicated in the loss of *ccrAB*-mediated SCC element excision.

Alongside this general picture of gene family stability and proliferation, however, individual genes may also be lost. One interesting case is that of *sel29p*, observed in two ST890 strains isolated in or prior to the 1940s but not seen subsequently in any public sequence data set. Could this gene have become extinct during the last 70 years? That it was seen only in a pseudogenized form could lend weight to such a hypothesis, as pseudogenization may lead to gene excision or deterioration, particularly during host adaptation (31). To put the absence of copies of *sel29*/*sel29p* within the GenBank nr/nt and RefSeq databases into sharper view, we attempted to determine the sequence types of all 11,351 *S. aureus* genomes downloaded from RefSeq and were able to easily achieve unambiguous predictions without manual intervention in 11,026 cases (97%; see Data Set S2 in the supplemental material for all predictions). We failed to find any further ST890 strains within this data set. Consequently, an alternative hypothesis of restriction of this gene to the ST890 lineage cannot be ruled out. As well as genes that are widespread across *S. aureus* strains, we have seen cases of genes restricted to a narrow range of lineages (e.g., *sel33*), so this scenario would not



be unprecedented. We further note that the lineage itself has not become extinct, with recent reports of ST890 strains derived from small mammals (32, 33). However, it does not appear that any of these strains have yet been subjected to whole-genome sequencing. It will be interesting to discover patterns of *sel29/sel29p* presence and absence within these strains should sequencing data become available, particularly as at least one of the two NCTC ST890 strains was isolated from a different host (human; see Data Set S1).

Our analyses indicated an association between clonal complex and the number of SE genes/gene clusters in the strains investigated. It would be interesting to further investigate possible associations between CCs and SE gene profiles (e.g., the particular pattern of SE genes that a strain possesses), though the larger RefSeq data set may be required to achieve statistical significance. Many studies have indicated associations between SE gene profiles and disease type, for example, between the *egc* gene cluster and both cystic fibrosis (34) and toxic shock syndrome (35). Other studies have shown links between disease and both SE profile and CC, such as those between CC30, infective endocarditis, and the genes *tsst-1*, *sea*, *sed*, *see*, and *sei* (36). Differences in regulatory system (7) may contribute to disease/SE gene associations, with some SE genes more likely to occur in chronic rather than acute infections. The plasticity of the mobile genetic elements carrying the majority of SE genes, with the different variant combinations circulating (e.g., see Fig. S2b), and the potential for their rapid loss and gain mean that selection could act swiftly on SE gene profiles. Considering such a system, it would seem plausible that strains with particular SE gene profiles would be selected for their roles in specific diseases and that clonal expansion would subsequently drive (at least some) CCs specialized for certain diseases. Limited within-CC recombination (37) would preserve these associations, establishing the patterns we see among extant and preserved strains. However, SE gene colocation might also mean some observed associations are indirect. A substantial meta-analysis of sequenced strains with high-quality disease status would undoubtedly be illuminating. It would be interesting to investigate, for example, whether the absence of *selx* on CC30 strains and resulting negative associations between this gene and other genes prevalent (e.g., *tsst-1* and *selu*) in this CC are due to simple gene loss and subsequent clonal expansion or due to strong selection for gene content.

Possessing multiple SEs may also be an advantage in itself. Distinct SE constituents of the *egc* gene cluster have been shown to exhibit different  $V\beta$  specificities and are therefore likely to have complementary effects on a host's immune system (22). Strains possessing multiple SEs could therefore possess a selective advantage regarding host colonization/invasion. Additionally, should two or more SEs be genetically linked, such as in the *egc* gene cluster, there are further opportunities for production of novel SE forms through processes such as recombination, such as we have seen here with the single observed case to date of *sel33*. Indeed, the high prevalence of distinct, nontrivial SE combinations presents a problem to researchers attempting to produce SE toxin-based vaccines for *S. aureus* infections. Despite promising results concerning strains with simple SE profiles (38), producing such a vaccine against strains with multiple SE genes remains a challenge (39). Consequently, disease-specific SE-based vaccines may be required, with a tailored combination of antitoxins.

Of the 845 SE and *tsst-1* gene sequences identified in this study, 123 (14.6%) were potential or likely pseudogenes. Approximately half of these cases were the *sel26* gene. The high rate of pseudogenization of this gene is potentially a consequence of its chromosomal location, as excision is not so easily possible as it is for the majority of SE genes residing on mobile elements. That said, we see a much lower rate of pseudogenization for *selx*, another chromosomal gene with only 4.6% cases present in strains with distinct sequence types. The different rates of the two genes may be due to the age of the genes or their importance in certain environmental (e.g., disease) niches, which further research might uncover. In general, rates of pseudogenization vary both between genes and between strains, indicating that selection may have played a key

role in this process. Eleven strains possess three or more pseudogenes. For example, NCTC 6133 has seven cases, including four of its six *egc* genes and the *selx* gene. Within the *egc* gene cluster, *seg* has the greatest number of pseudogenes. The observation of OMIUN *egc* arrangements in other strains may indicate that *seg* is not always essential to the success of the gene cluster. Most cases of pseudogenization do appear to be minor sequence changes to an established gene, which are likely to render it dysfunctional. However, the *ses-2p* gene adjacent to the *orfX* locus is more intriguing and has potential to be an emerging SE gene that has yet to be functional. Further sequencing of past and future *S. aureus* strains may shed light on the evolutionary trajectory of this sequence.

Here, we have shown how the analysis of even a medium-sized strain set can provide valuable information to the study of an important bacterial gene family. An added dimension to our analysis is that the strain set was collected over a period of almost a century, thereby granting access to biological material that can no longer be collected today. However, while the set is unique in terms of the specific strains involved, there are other collections worldwide that will now be contemplating or even carrying out sequence programs such as that which enabled this study. It will be interesting to see what information emerges.

## MATERIALS AND METHODS

**Data set preparation.** The genome assemblies of 133 *Staphylococcus aureus* strains (see Data Set S1 in the supplemental material for strain identities) derived from PacBio raw reads at the Wellcome Sanger Institute (WTSI) were downloaded in FASTA format from GenBank (<https://www.ncbi.nlm.nih.gov/genbank/>). The UniProt database (<https://www.uniprot.org/>) was searched for protein and nucleotide sequences attributed to each of the 25 target genes (24 SEs plus *tsst-1*) from other *S. aureus* strains. In the two cases where no SE matches were found (*selv* and *selw*), sequences were instead obtained from GenBank. Sequences were divided into two sets, where set 1 consisted of *selx* and *tsst-1* and set 2 consisted of the remaining 23 SEs.

**Enterotoxin gene hunting.** The software tool HMMER (40) was used to build profile hidden Markov models (pHMMs) for each gene set. The two pHMMs were then used to search the 133 genome assemblies for target gene matches. Searches were made using stringent parameters to guarantee full-length, or close to full-length, matches (for set 1,  $E < 1 \times 10^{-10}$ ,  $a < 5$ , and  $b > 600$ ; for set 2,  $E < 1 \times 10^{-10}$ ,  $a < 88$ , and  $b > 590$ ; where  $a$  is the starting coordinate of the match relative to the pHMM and  $b$  is the end coordinate of the match) and more relaxed parameters ( $E < 1 \times 10^{-10}$  for sets 1 and 2). Coordinates were chosen following visual inspection of the initial HMMER output for a nontrivial subset of strains.

The HMMER accessory tool Easel was subsequently used to extract the nucleotide sequences of all predicted target genes in all strains, keeping the two gene sets distinct. Gene-specific pHMMs were also built for each of the 25 target genes, and HMMER was again used to search for and extract predicted sequences, this time separated by gene identity. The gene-specific data sets were compared to the set-specific results for all strains to confirm that the gene family wide-approach was consistent with the gene-specific approach. No inconsistencies were identified. Extracted nucleotide sequences were aligned, along with 25 reference sequences (one for each target gene—see Table S1 in the supplemental material for the GenBank accession numbers of all SE gene references), using MUSCLE (41) and were manually divided into gene-specific groups within BioEdit (42). Using the reference sequences as guides, gene coordinates within the HMMER output were manually adjusted to ensure all sequences were full length. Using the modified coordinates, full-length target gene matches were extracted from the 133 *S. aureus* strains with Easel and realigned into gene-specific groups.

We also attempted both to gain support for the existence of the putative novel enterotoxin genes identified via this approach and to glean information on their origins by analyzing additional *S. aureus* genomes. Briefly, all 11,351 genomes within the RefSeq database (43) available on 11 May 2020 were downloaded and were subjected to an almost identical HMM-searching procedure as the 133 NCTC genome sequences. The only difference in the two procedures was the use of MAFFT (44) for gene sequence alignment, in place of MUSCLE, due to its ability to align tens of thousands of gene sequences within a few hours (using parameters *-retree 1 -maxiterate 0 -reorder*).

**Phylogenetic analysis.** Translated amino acid sequences of novel or recently identified gene-specific groups not included in the gene hunting process, one from each group, were aligned using MUSCLE along with reference sequences for previously known groups and the alignment input to the IQ-TREE phylogenetic software (27) with amino acid substitution model selection requested. During this process, the translated reference sequences for four groups (Gr5, Gr7, Gr11, and Gr12) were “repaired” with minor manual editing as the group members appeared to be pseudogenes with various mutations such as indel-causing frameshifts and premature stop codons. The repairs, which effectively estimated the amino acid states of the sequences before their putative pseudogenization but after their divergence from the other enterotoxin sequences, were made to maximize the phylogenetic signal in the data set and hence the reliability of the resulting tree. The resulting sequences are shown in Table S3. We also compared the nucleotide reference sequences of all established and putative gene family groups with the GenBank nr/nt database using BLAST (45).

**Strain typing.** Multilocus sequence typing (MLST) was conducted for each NCTC strain using the established seven-gene set for *S. aureus* (*arcC*, *aroE*, *glpF*, *gmk*, *pta*, *tpi*, and *yqjL*) (46). For each gene, all *S. aureus* sequences in the relevant BIGSdb database (47) were downloaded and a pHMM was calculated using HMMER. For each strain, the MLST gene sequences were extracted and concatenated into a single file, with the file subsequently input to BIGSdb for MLST characterization and identification of clonal complex. We also carried out spa typing for each strain, whereby the combination of differing repeat sequence types within the *SpA* gene was established. This process was carried out using the *get\_spa\_type.py* software ([https://github.com/mjsull/spa\\_typing](https://github.com/mjsull/spa_typing)), which compares repeats found between pairs of primer sequences against the Ridom (<https://spa.ridom.de/spatypes.shtml>) and eGenomics typing nomenclature. For 12 strains, manual editing of their genome sequence was required to achieve a spa type, as sequence mutations within their *SpA* genes meant that 100% matches to primer sequences were no longer achievable, preventing the software from extracting the repeats. All sequence type, clonal complex, and spa type predictions are shown in Data Set S1. The sequence type/clonal complex designation process was repeated for the 11,351 *S. aureus* genome assemblies downloaded from the RefSeq database. Predictions for the 11,026 strains (97%) for which manual intervention was not required to achieve a result are shown in Data Set S2.

**SNP analysis.** Each NCTC genome assembly was compared to that of NCTC 1803 (one full-length chromosome only, represented by a single contig) with Parsnp from the HarvestTools suite (26). The MFA file output from the suite's Gingr tool, which consisted of core genome single nucleotide polymorphisms across the 133-strain set, was used as input to the IQ-TREE phylogenetic analysis tool. Gingr was also used to visualize patterns of recombination between the strains.

**Statistical analysis.** The gene contents of strains for which a clonal complex origin was determined were analyzed along with year of isolation (set to the most recent year possible given the strain metadata shown in Data Set 1) within the R statistical environment (version 3.6.1) (48). A generalized linear model with a logarithmic link function and Poisson error distribution was fitted to the remaining data, with the number of genes/gene clusters (gene clusters were counted as if they were a single gene to account for the dependence of gene number counts on gene cluster presence and absence) as the independent variable and clonal complex (a factor with six levels) and year (after 1924) as dependent variables.

**Data availability.** All NCTC genomes have been deposited in the NCBI Sequence Read Archive under NCBI BioProject accession number [PRJEB6403](https://www.ncbi.nlm.nih.gov/bioproject/PRJEB6403). For individual accession numbers, please see Data Set S1, worksheet 1, in the supplemental material.

## SUPPLEMENTAL MATERIAL

Supplemental material is available online only.

**SUPPLEMENTAL FILE 1**, XLSX file, 5.1 MB.

**SUPPLEMENTAL FILE 2**, XLSX file, 0.6 MB.

**SUPPLEMENTAL FILE 3**, XLSX file, 0.02 MB.

**SUPPLEMENTAL FILE 4**, PDF file, 1.7 MB.

## ACKNOWLEDGMENTS

We are indebted to the work of all past NCTC staff and collaborators for their efforts in establishing, curating, and maintaining this collection, which at the time of writing is celebrating its 101st year of operation. We also express our gratitude to all who deposited strains into the NCTC and which were examined in this paper. We also thank three anonymous reviewers for their supportive and insightful comments.

The PacBio sequencing and genome assembly of the 133 *S. aureus* strains analyzed within this article were funded as part of the Wellcome Trust grant no. 101503/Z/13/Z "Creation of an e-resource centre to underpin the provision and use of type and reference strains of human pathogens."

## REFERENCES

- Graham PL, III, Lin SX, Larson EL. 2006. A U.S. population-based survey of *Staphylococcus aureus* colonization. *Ann Intern Med* 144:318–325. <https://doi.org/10.7326/0003-4819-144-5-200603070-00006>.
- Chowdhary VR, Tilahun AY, Clark CR, Grande JP, Rajagopalan G. 2012. Chronic exposure to staphylococcal superantigen elicits a systemic inflammatory disease mimicking lupus. *J Immunol* 189:2054–2062. <https://doi.org/10.4049/jimmunol.1201097>.
- McCormick JK, Yarwood JM, Schlievert PM. 2001. Toxic shock syndrome and bacterial superantigens: an update. *Annu Rev Microbiol* 55:77–104. <https://doi.org/10.1146/annurev.micro.55.1.77>.
- McCormick JK, Tripp TJ, Llera AS, Sundberg EJ, Dinges MM, Mariuzza RA, Schlievert PM. 2003. Functional analysis of the TCR binding domain of toxic shock syndrome toxin-1 predicts further diversity in MHC class II/superantigen/TCR ternary complexes. *J Immunol* 171:1385–1392. <https://doi.org/10.4049/jimmunol.171.3.1385>.
- Schlievert PM, Shands KN, Dan BB, Schmid GP, Nishimura RD. 1981. Identification and characterization of an exotoxin from *Staphylococcus aureus* associated with toxic-shock syndrome. *J Infect Dis* 143:509–516. <https://doi.org/10.1093/infdis/143.4.509>.
- Bergdoll MS, Crass BA, Reiser RF, Robbins RN, Davis JP. 1981. A new staphylococcal enterotoxin, enterotoxin F, associated with toxic-shock-syndrome *Staphylococcus aureus* isolates. *Lancet* i:1017–1021. [https://doi.org/10.1016/S0140-6736\(81\)92186-3](https://doi.org/10.1016/S0140-6736(81)92186-3).
- Fisher EL, Otto M, Cheung GYC. 2018. Basis of virulence in enterotoxin-mediated staphylococcal food poisoning. *Front Microbiol* 9:436. <https://doi.org/10.3389/fmicb.2018.00436>.

8. Langley RJ, Ting YT, Clow F, Young PG, Radcliff FJ, Choi JM, Sequeira RP, Holtfreter S, Baker H, Fraser JD. 2017. Staphylococcal enterotoxin-like X (SEIX) is a unique superantigen with functional features of two major families of staphylococcal virulence factors. *PLoS Pathog* 13:e1006549. <https://doi.org/10.1371/journal.ppat.1006549>.
9. Argudin MA, Mendoza MC, Rodicio MR. 2010. Food poisoning and *Staphylococcus aureus* enterotoxins. *Toxins (Basel)* 2:1751–1773. <https://doi.org/10.3390/toxins2071751>.
10. Thomas DY, Jarraud S, Lemerrier B, Cozon G, Echasserieau K, Etienne J, Gougeon ML, Lina G, Vandenesch F. 2006. Staphylococcal enterotoxin-like toxins U2 and V, two new staphylococcal superantigens arising from recombination within the enterotoxin gene cluster. *Infect Immun* 74:4724–4734. <https://doi.org/10.1128/IAI.00132-06>.
11. Okumura K, Shimomura Y, Murayama SY, Yagi J, Ubukata K, Kirikae T, Miyoshi-Akiyama T. 2012. Evolutionary paths of streptococcal and staphylococcal superantigens. *BMC Genomics* 13:404. <https://doi.org/10.1186/1471-2164-13-404>.
12. Wilson GJ, Tuffs SW, Wee BA, Seo KS, Park N, Connelley T, Guinane CM, Morrison WI, Fitzgerald JR. 2018. Bovine *Staphylococcus aureus* superantigens stimulate the entire T cell repertoire of cattle. *Infect Immun* 86:e00505-18. <https://doi.org/10.1128/IAI.00505-18>.
13. Lina G, Bohach GA, Nair SP, Hiramatsu K, Jouvin-Marthe E, Mariuzza R, International Nomenclature Committee for Staphylococcal Superantigens. 2004. Standard nomenclature for the superantigens expressed by *Staphylococcus*. *J Infect Dis* 189:2334–2336. <https://doi.org/10.1086/420852>.
14. Aung MS, Urushibara N, Kawaguchiya M, Sumi A, Takahashi S, Ike M, Ito M, Habadera S, Kobayashi N. 2019. Molecular epidemiological characterization of *Staphylococcus argenteus* clinical isolates in Japan: identification of three clones (ST1223, ST2198, and ST2550) and a novel staphylocoagulase genotype XV. *Microorganisms* 7:389. <https://doi.org/10.3390/microorganisms7100389>.
15. IWG-SCC. 2009. Classification of staphylococcal cassette chromosome *mec* (SCCmec): guidelines for reporting novel SCCmec elements. *Antimicrob Agents Chemother* 53:4961–4967. <https://doi.org/10.1128/AAC.00579-09>.
16. Luong TT, Ouyang S, Bush K, Lee CY. 2002. Type 1 capsule genes of *Staphylococcus aureus* are carried in a staphylococcal cassette chromosome genetic element. *J Bacteriol* 184:3623–3629. <https://doi.org/10.1128/jb.184.13.3623-3629.2002>.
17. Boundy S, Safo MK, Wang L, Musayev FN, O'Farrell HC, Rife JP, Archer GL. 2013. Characterization of the *Staphylococcus aureus* rRNA methyltransferase encoded by orfX, the gene containing the staphylococcal chromosome cassette *mec* (SCCmec) insertion site. *J Biol Chem* 288:132–140. <https://doi.org/10.1074/jbc.M112.385138>.
18. Zhang D-F, Yang X-Y, Zhang J, Qin X, Huang X, Cui Y, Zhou M, Shi C, French NP, Shi X. 2018. Identification and characterization of two novel superantigens among *Staphylococcus aureus* complex. *Int J Med Microbiol* 308:438–446. <https://doi.org/10.1016/j.jimm.2018.03.002>.
19. Noto MJ, Kreiswirth BN, Monk AB, Archer GL. 2008. Gene acquisition at the insertion site for SCCmec, the genomic island conferring methicillin resistance in *Staphylococcus aureus*. *J Bacteriol* 190:1276–1283. <https://doi.org/10.1128/JB.01128-07>.
20. Tunsjø HS, Kalyanasundaram S, Worren MM, Leegaard TM, Moen AEF. 2017. High frequency of occupied attB regions in Norwegian *Staphylococcus aureus* isolates supports a two-step MRSA screening algorithm. *Eur J Clin Microbiol Infect Dis* 36:65–74. <https://doi.org/10.1007/s10096-016-2771-0>.
21. Munson SH, Tremaine MT, Betley MJ, Welch RA. 1998. Identification and characterization of staphylococcal enterotoxin types G and I from *Staphylococcus aureus*. *Infect Immun* 66:3337–3348. <https://doi.org/10.1128/IAI.66.7.3337-3348.1998>.
22. Jarraud S, Peyrat MA, Lim A, Tristan A, Bes M, Mougé C, Etienne J, Vandenesch F, Bonneville M, Lina G. 2001. egc, a highly prevalent operon of enterotoxin gene, forms a putative nursery of superantigens in *Staphylococcus aureus*. *J Immunol* 166:669–677. <https://doi.org/10.4049/jimmunol.166.1.669>.
23. Letertre C, Perelle S, Dilasser F, Fach P. 2003. Identification of a new putative enterotoxin SEU encoded by the egc cluster of *Staphylococcus aureus*. *J Appl Microbiol* 95:38–43. <https://doi.org/10.1046/j.1365-2672.2003.01957.x>.
24. Varshney AK, Mediavilla JR, Robiou N, Guh A, Wang X, Gialanella P, Levi MH, Kreiswirth BN, Fries BC. 2009. Diverse enterotoxin gene profiles among clonal complexes of *Staphylococcus aureus* isolates from the Bronx, New York. *Appl Environ Microbiol* 75:6839–6849. <https://doi.org/10.1128/AEM.00272-09>.
25. Baba T, Bae T, Schneewind O, Takeuchi F, Hiramatsu K. 2008. Genome sequence of *Staphylococcus aureus* strain Newman and comparative analysis of staphylococcal genomes: polymorphism and evolution of two major pathogenicity islands. *J Bacteriol* 190:300–310. <https://doi.org/10.1128/JB.01000-07>.
26. Treangen TJ, Ondov BD, Koren S, Phillippy AM. 2014. The Harvest suite for rapid core-genome alignment and visualization of thousands of intraspecific microbial genomes. *Genome Biol* 15:524. <https://doi.org/10.1186/s13059-014-0524-x>.
27. Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol* 32:268–274. <https://doi.org/10.1093/molbev/msu300>.
28. Pushker R, Mira A, Rodríguez-Valera F. 2004. Comparative genomics of gene-family size in closely related bacteria. *Genome Biol* 5:R27. <https://doi.org/10.1186/gb-2004-5-4-r27>.
29. Sanchez-Herrero JF, Bernabeu M, Prieto A, Hüttner M, Juárez A. 2020. Gene duplications in the genomes of staphylococci and enterococci. *Front Mol Biosci* 7:160. <https://doi.org/10.3389/fmolb.2020.00160>.
30. Moon BY, Park JY, Hwang SY, Robinson DA, Thomas JC, Fitzgerald JR, Park YH, Seo KS. 2015. Phage-mediated horizontal transfer of a *Staphylococcus aureus* virulence-associated genomic island. *Sci Rep* 5:9784. <https://doi.org/10.1038/srep09784>.
31. Goodhead I, Darby AC. 2015. Taking the pseudo out of pseudogenes. *Curr Opin Microbiol* 23:102–109. <https://doi.org/10.1016/j.mib.2014.11.012>.
32. Monecke S, Gavriel-Widén D, Hotzel H, Peters M, Guenther S, Lazaris A, Loncaric I, Müller E, Reissig A, Ruppelt-Lorz A, Shore AC, Walter B, Coleman DC, Ehrlich R. 2016. Diversity of *Staphylococcus aureus* isolates in European wildlife. *PLoS One* 11:e0168433. <https://doi.org/10.1371/journal.pone.0168433>.
33. Mrochen DM, Schulz D, Fischer S, Jeske K, El Gohary H, Reil D, Imholt C, Trübe P, Suchomel J, Tricaud T, Jacob J, Heroldová M, Bröker BM, Strommenger B, Walther B, Ulrich RG, Holtfreter S. 2018. Wild rodents and shrews are natural hosts of *Staphylococcus aureus*. *Int J Med Microbiol* 308:590–597. <https://doi.org/10.1016/j.jimm.2017.09.014>.
34. Fischer AJ, Kilgore SH, Singh SB, Allen PD, Hansen AR, Limoli DH, Schlievert PM. 2019. High prevalence of *Staphylococcus aureus* enterotoxin gene cluster superantigens in cystic fibrosis clinical isolates. *Genes* 10:1036. <https://doi.org/10.3390/genes10121036>.
35. Jarraud S, Cozon G, Vandenesch F, Bes M, Etienne J, Lina G. 1999. Involvement of enterotoxins G and I in staphylococcal toxic shock syndrome and staphylococcal scarlet fever. *J Clin Microbiol* 37:2446–2449. <https://doi.org/10.1128/JCM.37.8.2446-2449.1999>.
36. Nienaber JJ, Sharma Kuinkel BK, Clarke-Pearson M, Lamlertthong S, Park L, Rude TH, Barriere S, Woods CW, Chu VH, Marín M, Bukovski S, Garcia P, Corey GR, Korman T, Doco-Lecompte T, Murdoch DR, Reller LB, Fowler VG, Jr, International Collaboration on Endocarditis-Microbiology Investigators. 2011. Methicillin-susceptible *Staphylococcus aureus* endocarditis isolates are associated with clonal complex 30 genotype and a distinct repertoire of enterotoxins and adhesins. *J Infect Dis* 204:704–713. <https://doi.org/10.1093/infdis/jir389>.
37. Driebe EM, Sahl JW, Roe C, Bowers JR, Schupp JM, Gillette JD, Kelley E, Price LB, Pearson TR, Hepp CM, Brzoska PM, Cummings CA, Furtado MR, Andersen PS, Stegger M, Engelthaler DM, Keim PS. 2015. Using whole genome analysis to examine recombination across diverse sequence types of *Staphylococcus aureus*. *PLoS One* 10:e0130955. <https://doi.org/10.1371/journal.pone.0130955>.
38. Aman MJ. 2017. Superantigens of a superbug: major culprits of *Staphylococcus aureus* disease? *Virulence* 8:607–610. <https://doi.org/10.1080/21505594.2016.1255399>.
39. Aguilar JL, Varshney AK, Pechuan X, Dutta K, Nosanchuk JD, Fries BC. 2017. Monoclonal antibodies protect from staphylococcal enterotoxin K (SEK) induced toxic shock and sepsis by USA300 *Staphylococcus aureus*. *Virulence* 8:741–750. <https://doi.org/10.1080/21505594.2016.1231295>.
40. Wheeler TJ, Eddy SR. 2013. nhmmer: DNA homology search with profile HMMs. *Bioinformatics* 29:2487–2489. <https://doi.org/10.1093/bioinformatics/btt403>.
41. Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32:1792–1797. <https://doi.org/10.1093/nar/gkh340>.
42. Hall TA. 1999. BioEdit: a user-friendly biological sequence alignment editor and analysis program for windows 95/98/NT. *Nucleic Acids Symp Ser* 41:95–98.
43. O'Leary N, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badredin A, Bao Y, Blinkova O, Brover V, Chetverin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D, Hatcher E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey KM, Murphy MR, O'Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W, Landrum MJ, Kimchi A, Tatusova T, DiCuccio M, Kitts P, Murphy TD, Pruitt KD. 2016. Reference sequence (RefSeq) database at NCBI: current

- status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* 44:D733–D745. <https://doi.org/10.1093/nar/gkv1189>.
44. Katoh K, Standley DM. 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol* 30:772–780. <https://doi.org/10.1093/molbev/mst010>.
45. 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](https://doi.org/10.1016/S0022-2836(05)80360-2).
46. Enright MC, Day NPJ, Davies CE, Peacock SJ, Spratt BG. 2000. Multilocus sequence typing for characterization of methicillin-resistant and methicillin-susceptible clones of *Staphylococcus aureus*. *J Clin Microbiol* 38:1008–1015. <https://doi.org/10.1128/JCM.38.3.1008-1015.2000>.
47. Jolley KA, Maiden MCJ. 2010. BIGSdb: scalable analysis of bacterial genome variation at the population level. *BMC Bioinformatics* 11:595. <https://doi.org/10.1186/1471-2105-11-595>.
48. R Core Team. 2019. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. <https://www.R-project.org>.
49. Letunic I, Bork P. 2007. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* 23:127–128. <https://doi.org/10.1093/bioinformatics/btl529>.