

1 Genome sequencing of an historic 2 *Staphylococcus aureus* collection reveals 3 new enterotoxin genes and sheds light on 4 the evolution and genomic organisation 5 of this key virulence gene family

6 **Jo Dicks,^a Jake Turnbull,^a Julie Russell,^a Julian Parkhill,^b Sarah Alexander,^{a,*}**
7 Culture Collections, Public Health England, London, United Kingdom^a; Department of Veterinary Medicine,
8 University of Cambridge, Cambridge, CB3 0ES, United Kingdom^b

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

26 **IMPORTANCE** *Staphylococcus aureus* strains pose a significant health risk to both
27 human and animal populations. Key amongst this species' virulence factors are the
28 Staphylococcal enterotoxin gene family. Certain enterotoxin forms can induce a po-
29 tentially life-threatening immune response, while others are implicated in less fatal
30 though often severe conditions such as food poisoning. Genetic characterisation of
31 Staphylococcal enterotoxin gene family members has steadily accumulated over re-
32 cent decades, with over 20 genes now established in the literature. Despite the cur-
33 rent wealth of knowledge on this important gene family, questions remain about the
34 presence of additional enterotoxin genes and the genomic composition of family mem-
35 bers. This study further expands knowledge of the Staphylococcal enterotoxins while
36 shedding light on their evolution over the last century.

37 **KEYWORDS:** *Staphylococcus aureus*, enterotoxin gene family, genome analysis,
38 National Collection of Type Cultures.

Compiled February 4, 2021

This is a draft manuscript, pre-submission

Address correspondence to Sarah Alexander,
Sarah.Alexander@phe.gov.uk.

39 INTRODUCTION

40 *Staphylococcus aureus* is a Gram-positive, coccoid bacterium belonging to the Firmi-
 41 cutes phylum of mainly low G+C bacteria. *S. aureus* is a common member of the hu-
 42 man microbiota, with studies estimating approximately 20-30% of the population to
 43 be long-term carriers of *S. aureus* strains in the skin, nostrils or female lower repro-
 44 ductive tract (1). In addition to its prevalence as a commensal organism of humans
 45 and animals, *S. aureus* is an important opportunistic pathogen. Strains can produce
 46 a variety of exotoxins, key amongst which are the staphylococcal enterotoxins (SEs),
 47 emetic toxins widely implicated in food poisoning. Gene family members are also as-
 48 sociated with more severe, life-threatening conditions. For example, SEB is classified
 49 as a potential bioterrorism threat given its rapid and acute stimulation of the immune
 50 system, and it is also potentially implicated in the inducement of auto-immunity (2).
 51 Toxic shock syndrome (TSS) is a serious, and potentially fatal, condition with roughly
 52 half of cases denoted as menstrual-associated and the remainder as non menstrual-
 53 associated. TSST-1, a protein very closely related to the SEs, gives rise to the majority
 54 of menstrual-associated TSS cases and approximately half of the non menstrual cases,
 55 with the remainder - ~25% in total - associated with SEB and SEC (3, 4).

56 The SE and TSST-1 proteins are superantigens (SAGs), immunomodulatory toxins
 57 that have the ability to stimulate large populations of T cells by interacting with the
 58 Variable region of the β -chain ($V\beta$) of the T-cell receptor. Structurally, SAGs are two-
 59 domain proteins characterised by a β -grasp domain and an OB-fold domain. The SE
 60 proteins are encoded by a family of genes related by their DNA sequence. The recent
 61 literature on the Staphylococcal enterotoxin gene family encompasses 24 genes: *sea*,
 62 *seb*, *sec*, *sed*, *see*, *seg*, *seh*, *sei*, *selj*, *sek*, *sel*, *sem*, *sen*, *seo*, *sep*, *seq*, *ser*, *ses*, *set*, *selu*, *selv*,
 63 *selw* (formerly *selu2*), *selx* and *sely*. The protein encoded by the closely related *tsst*-
 64 *1* gene was initially discovered independently by two groups as PEC (staphylococcal
 65 pyrogenic exotoxin C; 5) and SEF (staphylococcal enterotoxin F; 6) and later renamed
 66 TSST-1 upon agreement, hence the absence of the *sef* nomenclature within the SE
 67 gene list. Only genes whose proteins have demonstrated emetic activity are given the
 68 "se" prefix, with others designated as "sel" for "staphylococcal enterotoxin-like". The
 69 nomenclature used here is largely taken from Fisher et al. (7).

70 In phylogenetic terms, the majority of the SE genes group closely to one another
 71 and with a number of *Streptococcus pyogenes* genes (8). The *selx* and *tsst-1* gene se-
 72 quences group more distantly, with those of approximately 26 Staphylococcal Super-
 73 antigen-like exoproteins (SSLs), which unlike the SAGs are immune invasion molecules
 74 and will not be considered here further. Despite its phylogenetic placement amongst
 75 the SSLs, *selx* is functionally similar to the SAG genes and is therefore referred to as an
 76 "SSL-like SAG" (8). Uniquely, SEIX is a single-domain SAG, lacking the OB-fold domain
 77 seen in all other staphylococcal SAGs to date.

78 The majority of the SAG genes are located on mobile genetic elements such as
 79 pathogenicity islands, prophages and plasmids (7, 9). Clustering of the enterotoxin
 80 genes is also observed, most notably the *egc* cluster, which in a given strain can com-
 81 prise up to seven genes and pseudogenes from a repertoire of nine (pseudo)gene
 82 forms (10). Consequently, while there is considerable variability with regard to the en-
 83 terotoxin gene content between strains, the co-occurrence of the individual genes is
 84 highly non-random.

85 The National Collection of Type Cultures (NCTC) was founded in 1920 to address a
 86 recognised need for accumulating and disseminating information on human, animal,
 87 fungal and plant pathogens. It is one of four culture collections operated by Public
 88 Health England as part of a globally recognised biological resource centre, providing

89 many thousands of historical and emerging strains to researchers and biomedical sci-
90 entists worldwide. Recently, a Wellcome-funded initiative to sequence the genomes
91 of ~ 3,000 NCTC strains was completed. Amongst the datasets developed were Pa-
92 cific Bioscience (PacBio) long-read sequences and associated genome assemblies for
93 133 *Staphylococcus aureus* strains accessioned between 1924 and 2016, with at least
94 one strain isolated prior to June 1924. NCTC strains are accessioned either proactively,
95 based upon scientific requests from one of its nine past and current curators or pas-
96 sively, deposited by members of the research community. While we think it unlikely
97 that such a moderately-sized dataset would be representative - either geographically
98 or temporally - of globally circulating *S. aureus* strains over the last century, the dataset
99 is nonetheless diverse, with 43 distinct Sequence Types each represented by one or
100 more strains.

101 Here, we analysed this new dataset on an historic strain collection to answer ques-
102 tions on the number of staphylococcal enterotoxin genes and their genomic organ-
103 isation. We showed that each examined strain possessed between two and eigh-
104 teen SE genes. We identified seven putative SE genes outside our search list, four
105 of which were not seen in NCTC strains accessioned after 1951, and one of which is
106 the most prevalent enterotoxin-like sequence identified to date. We also examined
107 the genomes of over 11,000 *Staphylococcus aureus* strains in the RefSeq database in
108 order to gain support for this expanded SE gene repertoire. The RefSeq dataset of-
109 fered significant support for the newly identified genes while additionally presenting
110 a further three SE genes and two gene variants. Collectively the two datasets shed
111 light on the genomic distribution of SE genes, further delineating six gene clusters
112 and introducing a new one. Crucially, the NCTC dataset enabled the examination of
113 temporal patterns of enterotoxin birth and death to be made over a period of a cen-
114 tury, showing a remarkable stability of gene content over this time, with all but one
115 gene well represented in global sequence datasets. Finally, in accordance with this ob-
116 served stability, analysing the inter-relationships between the NCTC strains showed
117 that their Clonal Complex origins were more important than their time of isolation in
118 determining their enterotoxin gene load.

119 RESULTS

120 **The NCTC strains revealed novel enterotoxin-like sequences.** Using a profile
121 Hidden Markov Model approach to hunt for DNA sequences within a set of genomes
122 provides the opportunity to find novel SE genes that have yet to be formally charac-
123 terised. We identified 825 SE- and 20 *tsst-1*-like sequences within the 133 *S. aureus*
124 strains, with genomic co-ordinates and annotation details provided within Data Set
125 S1. The 845 sequences clustered into 29 easily distinguishable gene-specific groups.
126 In addition to finding 22 of the 25 expected SE/*tsst-1* gene-specific groups (reference
127 sequences for each group are shown in Table S1; no copies of *see*, *ses* or *set* were iden-
128 tified in any of the 133 strains), seven additional putative SE genes were identified,
129 which we initially termed Gr1-Gr7. Strikingly, 133 copies of the Gr1 sequence were
130 found, one in every strain analysed, spanning 42 distinct Sequence Types (including
131 three not found in the *S. aureus* BIGSdb). While almost half (62 out of 133 copies) were
132 likely to be pseudogenes due to premature stop codons or frameshift-inducing indels,
133 seemingly intact versions of the coding sequence were seen across the accessioning
134 period, with the most recent confirmed copy seen in strain NCTC 13434, isolated and
135 accessioned in 2008.

136 All other groups (denoted Gr2 to Gr7) consisted of between two and sixteen mem-
137 bers. The sixteen copies of Gr2 were found in strains isolated between 1932 and 2008.

138 Gr7 was less prevalent, with only five copies, but was observed over only a slightly
 139 reduced timespan, between 1938 and 1997. Despite its wide timespan, all five copies
 140 of Gr7 are likely to be pseudogenes. The other four genes were limited to the earlier
 141 strains, with Gr3/Gr4 (co-located), Gr5 and Gr6 most recently seen in strains acces-
 142 sioned in 1951, 1949 and 1948 respectively. Furthermore, both copies of Gr5 were
 143 presumed pseudogenes, likely the result of two small deletions. However, the preva-
 144 lence of pseudogeny outside the cases of Gr1, Gr5 and Gr7 was generally less common,
 145 with fewer potential or likely pseudogenes in all other gene groups (see Data Set S1
 146 for details).

147 **Origins of the putative enterotoxin-like sequences.** We searched for close se-
 148 quence matches to each of the seven initially unidentified enterotoxin-like sequences
 149 in order to establish whether they had been observed previously. Overall, we found
 150 882 high-scoring *megablast* hits to the GenBank nr/nt database (on 12/02/2020), with
 151 sequence-specific frequencies shown in Table S2a.

- 152 • **Gr1 is the near-ubiquitous chromosomal gene *sel26*.** The majority of the
 153 *megablast* hits, numbering 735 (83%, excluding 32 hits to NCTC genome se-
 154 quences), were to the Gr1 nucleotide sequence at 93.24-100% sequence iden-
 155 tity, including two full-length copies in a single strain (*S. aureus* strain ch22 chro-
 156 mosome; CP017807.1) and all but 17 of which were to complete genomes or
 157 chromosomes. The 17 gene hits were all annotated as enterotoxin-like W genes
 158 (e.g. 98.14% sequence identity to a purported *selw* gene identified in strain
 159 TD101; KX655716.1), though this gene was significantly different from the *egc*
 160 gene cluster *selw* gene used in our HMM search. It would appear that two
 161 distinct genes have been using the *selw* nomenclature: the *egc* gene formerly
 162 known as *selu2* and the seemingly ubiquitous (or at least highly prevalent) chro-
 163 mosomal gene previously seen in studies of human (11) and bovine (12) SEs.
 164 We will henceforth refer to this chromosomal gene as *sel26*, generally in line
 165 with the recommended nomenclature for SE genes (13) and similar to that used
 166 within (12), though currently lacking experimental confirmation to the best of
 167 our knowledge.
- 168 • **Gr2 is the *orfX*-associated gene *selz*.** We found 40 hits to the Gr2 nucleotide
 169 sequence, all at 96.54-99.49% sequence identity. Fourteen hits were to gene se-
 170 quences, the majority to a gene recently described as a staphylococcal entero-
 171 toxin-like Z gene (*selz*) in strains of *Staphylococcus argenteus* (14). Interestingly,
 172 four hits (e.g. KT316803.1) were annotated as a Staphylococcal Cassette Chro-
 173 mosome *mec* (*SCCmec*) element, a mobile genetic element implicated in broad-
 174 spectrum beta-lactam resistance via the *mecA* gene (15). An additional hit (U10927.2)
 175 appears to lie adjacent to an *SCCcap1* element, an SCC element with structural
 176 similarities to *SCCmec* but which instead harbours a type 1 capsular polysac-
 177 charide biosynthesis gene cluster (16). Further investigation of the annotation
 178 of U10927.2 shows that *selz* is the enterotoxin gene observed but unnamed in
 179 2002 by Luong et al. (16). We searched for proximity of *selz* to *orfX* (which en-
 180 codes an RlmH-type ribosomal methyltransferase) within each of the 16 strains
 181 with copies of Gr2/*selz*, as SCC elements are known to insert within the C ter-
 182 minus of this locus (17). We found that all copies of Gr2 were indeed in close
 183 proximity to *orfX* (between 2.5kb and 38.6kb), with two strains (NCTC 10399 and
 184 NCTC 10649) possessing an adjacent *SCCcap* element but none an *SCCmec* el-
 185 element. Six Gr2/*selz* copies were found in strains with Sequence Type 121, with
 186 others belonging to ST123, ST151, ST351, ST395, ST705, ST707, ST1254 and a
 187 novel Sequence Type. None of these strains belong to the six major identified

188 Clonal Complex groups (CC1, CC5, CC8, CC22, CC30 and CC97) in this study.

189 • **Gr3 and Gr4 are the clustered genes *sel27* and *sel28*.** The clustered Gr3 and
 190 Gr4 nucleotide sequences found 39 hits each, up to 98.41% and 99.12% se-
 191 quence identity respectively. Of particular note were strong hits for each gene
 192 to a gene cluster pathogenicity island in strain 364P and to two recently iden-
 193 tified enterotoxin genes annotated as *Sel27* and *Sel28* in strains SJTU F20365
 194 (MF370878.1; 18), 86, 72, 50, SG19, SG16, SG13, SG11, SG09, SG05-2, SG05-1,
 195 SG04 and SG01. The oldest confirmed NCTC strain identified as carrying these
 196 two genes, NCTC 5664, was isolated in 1936 and the putative youngest, NCTC
 197 8765, during or prior to 1951. Five of the seven NCTC strains possessing the
 198 two genes belong to ST9 (CC1) and the remainder to ST350.

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

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

211 • **Gr7 is the *orfX*-associated pseudogene *ses-2p*, a variant of *ses* clustered
 212 with *seh*.** Finally, the Gr7 nucleotide sequence, present in all five NCTC copies in
 213 close proximity to an *seh* gene sequence, found 21 hits ranging between 93.51%
 214 and 100% sequence identity. While the majority were to complete genome or
 215 chromosome sequences, two hits (EU272079.1 and KX690110.1) were - similar
 216 to Gr2 - to insertion sites of SCCmec elements. Further investigation of these
 217 hits identified previous reports of a partial enterotoxin gene with sequence sim-
 218 ilarity to *seo*, in close proximity to an *seh* gene and associated with SCCmec
 219 Type IV element insertion (19, 20). Similar to our analysis of Gr2, we searched
 220 for *orfX* genes within the five strains with copies of both Gr7 and *seh*, plus the
 221 single strain (NCTC 13435) possessing a presumed pseudogenised version of
 222 *seh* but no copy of Gr7. We found that all five copies of Gr7 and six copies of
 223 *seh/seh_p* were close to an *orfX* gene (between 17.5kb and 43.1kb). However,
 224 only one copy of Gr7 (in NCTC 13297) was adjacent to an SCC element (in this
 225 case likely an SCC*fus* element) while the *seh_p* copy in NCTC 13435 was adjacent
 226 to an SCCmec Type IV element. Further sequence analysis showed that while
 227 the 3' region of the Gr7 gene (~300bp) was highly similar to the corresponding
 228 region of *ses* (rather than *seo*), the 5' region of approximately 30bp was ~40bp
 229 shorter and dissimilar at the sequence level, with the intervening region show-
 230 ing a moderate level of sequence similarity interrupted by several presumed
 231 mutations. All copies of Gr7 in NCTC strains therefore look to be (only partially)
 232 truncated pseudogenes, containing several premature stop codons, though it
 233 is uncertain whether Gr7 was ever a functional gene. Four of the NCTC strains
 234 with the Gr7 sequence are ST10, with the remaining strain ST1. We propose
 235 that this pseudogene is referred to henceforth as *ses-2p*.

236 **Most novel enterotoxin-like gene and pseudogene sequences are also ob-**
 237 **served within RefSeq genomes.** We investigated the numbers of high-scoring hits

238 for each of the seven novel or recently identified (pseudo)genes within 11,351 *Staphy-*
 239 *lococcus aureus* genome sequences within the RefSeq database (on 11/05/2020). Of
 240 the 64,281 total hits to staphylococcal enterotoxin-like sequences, 12,505 were to six
 241 of these seven genes. Only Gr5/*sel29p* failed to find hits in strains other than the two
 242 NCTC strains (which were present within the RefSeq database). Table S2a shows that
 243 the relative frequencies of the other six genes within the NCTC dataset are mirrored
 244 within the *nr/nt* and *RefSeq* databases. This suggests to a certain extent that the fre-
 245 quencies of SE genes within the NCTC dataset are indicative of their frequencies within
 246 larger datasets, notwithstanding the absence of three SE genes from our dataset. Al-
 247 though we did find instances of *see*, *ses* and *set* within the sizeable RefSeq dataset (2,
 248 39 and 36 copies, respectively), their low frequencies suggest these may be relatively
 249 rare genes, or certainly within strains whose genomes have been sequenced thus far.

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

- 258 • **Gr8 and Gr9 are the clustered genes *sel31* and *sel32*.** Fifteen copies of the Gr8
 259 sequence in the RefSeq genomes were to a likely functional (based on its amino
 260 acid translation) gene which we term *sel31*. Of the eighteen Gr9 sequences, fif-
 261 teen directly neighboured *sel31*. We term this new gene *sel32*, and hence delin-
 262 eate a new enterotoxin gene cluster. Fifteen of the eighteen RefSeq genomes
 263 could be classified as four Sequence Types (ST1, ST121, ST97, ST508) from three
 264 Clonal Complexes (CC1, CC45, CC97), see Data Set S3 for details of these strains.
 265 Each gene found three identical BLAST hits to the GenBank *nr/nt* database, from
 266 the same genomic sources, two of which were to plasmid genomes, highlighting
 267 the likely origin of the gene pair.
- 268 • **Gr10 is the egc cluster gene *sel33*, a recombinant of *selw* and *sen*.** A sin-
 269 gle copy of a new recombinant derivative of egc gene cluster genes *selw* and
 270 *sen* was identified in strain BSAC1477 from the BSAC Resistance Surveillance
 271 Project (GenBank accession NZ_FGMI01000018.1; [BSAC](#)), which we term *sel33*.
 272 This strain was isolated in or after 2001 and derives from CC22. Of the five se-
 273 quences investigated here, only *sel33* failed to find any similar sequences within
 274 the *nr/nt* database, suggesting this recombination to be a rare occurrence.
- 275 • **Gr11 is the variant of *seh*, *seh-2p*, strongly associated with SCCmec Type**
 276 **IV elements.** The RefSeq dataset contained thirty copies of Gr11, a pseudo-
 277 genised or truncated form of *seh* which we will refer to henceforth as *seh-2p*.
 278 The sequence of *seh-2p* likely possesses two single nucleotide deletions relative
 279 to the canonical form of *seh*. All but two classified genomes derive from ST80,
 280 with one each of ST4563 and a novel Sequence Type, both differing in a only sin-
 281 gle MLST allele from ST80 (both within the *glpF* gene). Notably, one of the ST80
 282 strains is NCTC 13435, which unlike the five copies of *seh* in the NCTC strains,
 283 lacked a neighbouring *ses-2p* sequence. Similarly, none of the 30 copies of *seh-*
 284 *2p* within the RefSeq dataset have a neighbouring enterotoxin-like sequence,
 285 and 29 show significant evidence of a neighbouring SCCmecIVc(2B) element,
 286 with 27 genomes showing all SCCmec genes spread over one to four contigs
 287 and two additional genomes showing partial SCCmec matches (both including

288 *mecA* presence). In contrast, 219 of the 224 copies of *seh* possessed an adjacent
 289 *ses-2p* or *ses-3p* (see below) sequence, with 29 and 190 copies respectively.

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

302 **The novel enterotoxin gene sequences are spread across the SE phylogeny.**

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

318 **There are at least seven Staphylococcal enterotoxin gene clusters.** SE genes
 319 are known to sometime co-locate with others, often on plasmids or pathogenicity is-
 320 lands. The most striking example of this phenomenon is the *egc* gene cluster. The
 321 full characterisation of this operon has taken place in a stepwise fashion. The initial
 322 discovery of the *seg* and *sei* genes (21) was followed by identification of *sem*, *sen*, and
 323 *seo* in the neighbouring genomic regions, along with evidence of their co-transcription,
 324 while two pseudogenes (ϕ ent1 and ϕ ent2) were found between *sei* and *sen* (22). A sixth
 325 gene, *selu*, thought to be the product of deletions within ϕ ent1 and ϕ ent2) was identi-
 326 fied later (23). Most recently, *selw* (formerly *selu2*) and *selv* were identified (10). While
 327 the nucleotide sequence of *selw* was highly similar to that of *selu*, the main difference
 328 being a 15bp deletion in the former compared to the latter, and thought to be the re-
 329 sult of a different mutation of the ϕ ent1 and ϕ ent2 sequences from that hypothesised
 330 in *selu*, *selv* was found to be the product of a recombination of *sem* and *sei*.

331 The *egc* gene cluster appears to be highly prevalent within *S. aureus* genomes
 332 (14, 24). In this study, 59 *egc* gene clusters were found in 58 of the 133 *S. aureus*
 333 strains (43.6%), all but one seemingly complete. One of these 58 strains (NCTC 7972)
 334 possessed two gene clusters, with one of them the only case observed in this strain
 335 set of a large gap between any two *egc* genes, 13.6kb between *seo* and *sem*. In NCTC
 336 11963, the only clear case of an incomplete *egc* gene cluster, we identified genes *seg*
 337 and *seo* on two separate genomic contigs, such that the omission of intervening genes

338 may be the result of an incomplete genome assembly rather than their absence from
 339 the genome. Further investigation of the raw sequencing reads overlapping this re-
 340 gion showed the central region of the *egc* gene cluster to suffer from a very low read
 341 coverage but nonetheless to offer support for the existence of a full gene cluster most
 342 likely of type OMIUNG (read 27628 runs from the middle of *seo* to the middle of *seg*
 343 and is highly similar to the corresponding sequence of NCTC 2669, albeit with a single
 344 large run of T's breaking the alignment, presumably a sequencing artefact).

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

358 Frequencies of the distinct *egc* gene cluster arrangements identified in the NCTC
 359 strains are given in Table 1, showing that OMIUNG (i.e. the gene order *seo-sem-sei-*
 360 *selu-sen-seg*) and its close variant OMIWNG (including 'minor' pseudogenes of all six
 361 genes) are the predominant gene cluster variants, seen in this strain set in a ratio of
 362 2.6:1. As well as three strains isolated and accessioned in the 1930s or 1940s possess-
 363 ing an OMIUN variant (i.e. apparent absence of the *seg* gene), we see a recent strain
 364 (NCTC 13373, accessioned in 2005 and equivalent to ATCC 43300, a clinical isolate from
 365 Kansas) with the OVUNG form, the potentially rare gene *selv* the result of a recombina-
 366 tion between *sem* and *sei*. A comparison of the sequence of *selv* in NCTC 13373 to the
 367 canonical form in strain A900624 (10) from the French National Reference Center for
 368 Staphylococci (see Fig S1) shows evidence for the recombination between *sem* and *sei*
 369 in NCTC 13373 having occurred slightly closer to the 5' end of the sequence, though
 370 both events clearly took place within a central sequence region highly similar between
 371 the two progenitor genes. That observation, together with a high number of single nu-
 372 cleotide differences between the NCTC 13373 and *selv* reference sequences plus the
 373 alternative OVWNG form of the A900624 gene cluster, indicates that the two genes
 374 likely arose from two distinct recombination events. Note that the OVWNG form and
 375 the OMI33G form (i.e. containing the novel *sel33* recombinant of *selw* and *sen*), which
 376 we discovered in strain BSAC1477, were not observed within the NCTC dataset.

377 Five additional gene clusters were observed within the NCTC strains, as shown in
 378 Table 2. All but one instance of the 122 gene clusters, the broken *egc* cluster men-
 379 tioned above, were found in intact gene cluster form. Interestingly, the *sel27-sel28*
 380 gene cluster, found within 7 strains accessioned in the 1930s-50s, was seen to be lo-
 381 cated close to the *egc* gene cluster (2 of the 3 OMIUN strains and 5 of the 14 OMIWNG
 382 strains). The distance of the *egc* cluster to the *sel27-sel28* gene cluster ranged between
 383 17,465 and 39,464bp, with an average of approximately 27.6 kb.

384 While the *sek-seq* and *seh-ses-2p* gene clusters are comprised of genes within the
 385 same clade in Fig 1 (the cyan and orange groups are often referred to as a single clade
 386 elsewhere), the remaining four clusters contain genes spanning two or even three
 387 clades. In particular, the two most prevalent *egc* gene cluster arrangements (OMI-

388 UNG and OMIWNG) which account for 54 of its 59 copies, possess two genes from
 389 three main SE clades (shaded yellow, cyan and purple in Fig 1). It has been speculated
 390 that this divergence of the *egc* cluster genes may indicate the cluster's role as the pro-
 391 genitor of the majority of SE genes in *Staphylococcus aureus* (22). The seventh gene
 392 cluster, *sel31-sel32*, identified in fifteen RefSeq genome sequences from three Clonal
 393 Complexes, was not observed within the NCTC strains.

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

402 We found ten cases of a strain possessing two copies of the same gene and a
 403 single case of three gene copies (see Data Set S1). Notable examples include NCTC
 404 7415, in which we found both three copies of *tsst-1* and two copies of the *sec-sel*
 405 cluster, and NCTC 7972, which harboured two intact *egc* gene clusters. In the former
 406 case, the likelihood of two of the three *tsst-1* copies and one of the *sec* copies being
 407 pseudogenes within a 20kbp region of a single contig may contribute to this finding.

408 Consistent with other studies such as Varshney et al. (24), individual strains were
 409 found to possess numerous *SE/tsst-1* genes, with a range of 2 to 18 genes per strain
 410 and a mean of 6.35 (median of 6). While this value is slightly higher than the average 5
 411 SE genes per strain seen in (24), that prior study had looked at fewer genes, 19 of the
 412 37 genes examined here, which may have led to the lower gene counts.

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

423 The six gene clusters identified in the previous section are easily apparent as ei-
 424 ther single or sets of large red circles. Additional positive and negative associations are
 425 also apparent. Notable positive associations include *tsst-1* with the *sec-sel* gene clus-
 426 ter, *seb* with the *sek-seq* gene cluster, *sely* with *selz*, and *selw* with the *sel27-sel28*
 427 gene cluster. The former two associations are previously noted and likely the products of
 428 co-presence on SaPI_{m1/n1} (or SaPI_{bov1}) and SaPI₃ pathogenicity islands respectively.
 429 Examination of the relevant gene co-ordinates shows that the associations are without
 430 exception underpinned by co-presence on the same contig, though not by co-location.
 431 The most compact examples are the 11 cases of *seb/sek-seq*, which are approximately
 432 11kb apart in all strains. However, to the best of our knowledge the latter two associa-
 433 tions have not been observed prior to this study. *sely* has previously been seen only on
 434 the chromosome, so its association with the *orfX*-associated *selz* gene could be down
 435 to chance alone. The *selw/sel27-sel28* association is likely to involve particular forms of
 436 the *vSaβ* genomic island, given the known presence of the *egc* cluster on this mobile
 437 element (25).

438 Notable negative associations include *selx* with the OMIUNG form of the *egc* gene
 439 cluster and the prophage-encoded *sea/sek-seq* gene cluster combination with OMI-
 440 UNG and *sely*. Although Varshney et al. (24) note an absence of *egc* gene clusters
 441 within *seb*⁺ strains, we observe a more complex pattern in this dataset. Examining
 442 the 17 *seb*⁺ strains identified here, we find that 7 of the 8 strains isolated in or prior
 443 to 1949 carried the *egc* gene cluster whereas none of the 9 strains believed to derive
 444 from the 1950s onwards were found to possess it. Furthermore, no strains with *sel27*
 445 and *sel28* harboured a *seb* gene. Taken together, these observations suggest different
 446 SE gene combinations have circulated within the *S. aureus* population, some of which
 447 were restricted to particular timeframes, or perhaps to different parts of the *S. aureus*
 448 population. Fig S2b shows a depiction of the associations observed in this study and/or
 449 described in Argudin et al. (9).

450 **Phylogenetic and temporal patterns of Staphylococcal enterotoxin genes.**

451 We estimated a phylogenetic tree of the 133 NCTC *S. aureus* strains using Harvest-
 452 Tools (26) and IQ-TREE (27), based on 96,541 core genome SNPs, and annotated it
 453 with SE/*tsst-1* gene content and Clonal Complex group. It is immediately clear from
 454 Fig 3 that certain SE gene combinations are restricted to particular clades within the
 455 tree. For example, only CC1, CC5 and CC22 strains harbour the OMIWNG form of the
 456 *egc* gene cluster in this dataset. Indeed, the observed patterns of gene content explain
 457 many of the associations between genes described in the previous sections. For ex-
 458 ample, we see that *selx* is completely absent from CC30 (purple strip), a monophyletic
 459 group in which the OMIUNG form of the *egc* gene cluster is highly prevalent, thereby
 460 explaining the strong negative association between *selx* and *selu*.

461 Most clades cover a broad timespan (e.g. confirmed isolation periods of at least
 462 1941-1997 for CC1; 1948-1988 for CC5; 1932-2003 for CC8; 1928-2003 for CC30) with a
 463 potential span of 1933-1985 for CC97 (though only presently confirmed up to 1954)
 464 and only CC22 represented here by a narrow group of strains (1990-2005). While
 465 Fig 3 suggests gene content to be highly correlated with Clonal Complex, the mean
 466 number of enterotoxin genes per strain was found to be higher in strains isolated
 467 within the 1920s-1940s (7.60) than in the 1950s-2010s (5.65), irrespective of Clonal
 468 Complex membership. We examined the relationship between the number of entero-
 469 toxin genes per strain with Clonal Complex and year of isolation by fitting a generalised
 470 linear model to the data, collapsing gene clusters to single observations as described in
 471 the Materials and Methods section. We found that for the 90 NCTC strains with Clonal
 472 Complex designations (see Fig S4 for plots of the data), the number of enterotoxin
 473 genes/gene clusters was strongly associated with Clonal Complex identity ($p < 0.05$
 474 for three of the five factors when compared to CC1) but neither with year of isolation
 475 ($p = 0.274$) nor the interaction between Clonal Complex and time ($p > 0.425$ for all
 476 factor interactions). This suggests that enterotoxin gene content within Clonal Com-
 477 plexes has remained stable across the century of strain isolation and that the putative
 478 temporal difference in gene content described above may be due to sampling effects,
 479 with a higher frequency of strains harbouring the *egc* gene cluster isolated during the
 480 earlier period (60% vs. 34%).

481 HarvestTools Gingr plots of the core genome SNP alleles alongside the phyloge-
 482 netic tree (see Fig S3 for an example, with NCTC 1803 as a reference genome) also
 483 indicate horizontal transmission between disparate *S. aureus* clades has taken place.
 484 Consequently, Staphylococcal enterotoxin gene content may also be influenced by
 485 horizontal processes in addition to clonal expansion. In future it would be interest-
 486 ing to analyse whether, for example, the two cases of the *sel27-sel28* gene cluster
 487 outside CC1 (uppermost green circles in Fig 3) were due to horizontal transfer of the

488 pathogenicity island on which they are located.

489 DISCUSSION

490 We analysed a strain set of 133 *Staphylococcus aureus* strains from the UK National Col-
491 lection of Type Cultures, with a particular goal of understanding the complement of
492 enterotoxin genes captured within, thereby further enhancing the utility of the strains
493 for the benefit of the research community. While we did not initially anticipate uncov-
494 ering any potential novel genes, particularly given the size of the dataset and the lack
495 of an enterotoxin-focussed strategy for its collection, the use of a pHMM-profile ap-
496 proach allowed us to identify new sequences that we hope will be investigated further
497 by researchers in this area. Given the relative ease with which we found putative en-
498 terotoxins first within the NCTC dataset, and subsequently within the RefSeq database,
499 this leads us to speculate as to whether there might yet be other enterotoxin genes
500 left for others to uncover. While the NCTC and RefSeq datasets encompass a size-
501 able proportion of the global diversity of *S. aureus* strains, with 43 and 468 distinct
502 Sequence Types represented respectively, they will not have captured the full range.
503 Consequently, as yet unidentified enterotoxin genes may still be present in strains
504 whose genomes are currently outside the reach of strain and sequence collections.

505 Our study has also added to the understanding of the genomic organisation of the
506 enterotoxin genes, particularly via gene clusters carried by mobile genetic elements.
507 Gene families harboured by bacterial genomes are presented with an array of strate-
508 gies that enable them to thrive and mobilise. The Staphylococcal enterotoxin gene
509 family has indeed shown it is capable of exploiting many of these routes, from use of
510 plasmids, prophages, pathogenicity islands and genomic islands in addition to stable
511 chromosomal inheritance. Despite the consequent stability of the SE genes over the
512 past century, in general large gene families appear to be rare in many bacteria. A study
513 of the sequenced genomes of species including *Escherichia coli*, *Streptomyces pyogenes*
514 and *Chlamydomypha pneumoniae* showed limited numbers of gene families of size 20 or
515 over, with only ~10 such gene families in *S. aureus* strains Mu50, MW2 and N315 (28).
516 A more recent study found greater variation in the number of gene families between
517 strains but again the number of duplicated genes was relatively limited, with a maxi-
518 mum of 190 duplicates for 84 gene families across 473 strains (29). Furthermore, the
519 majority of duplicated genes in the latter study were thought to have a phage origin.
520 The findings here are consistent with this observation.

521 The egc gene cluster has clearly been a key component in the expansion of the
522 SE gene family. Interestingly, the genomic island harbouring the egc gene cluster was
523 recently found, similarly to pathogenicity islands, to be capable of mobilisation due
524 to a temperate bacteriophage (30). The egc is the only SE gene cluster to date that
525 has been shown to have produced novel recombinant SE genes and from the grow-
526 ing number of components on offer, six distinct combinations were observed in the
527 strains analysed here. As mentioned earlier, the egc gene cluster has been mooted
528 as a putative SE nursery, whereby the observed genetic diversity has been generated
529 by the processes of tandem duplication and subsequent divergence (22). The level of
530 variation observed in this study would seem consistent with that view. Looking at the
531 seven gene clusters in Table 2, all but one (*sek - seq*) has members from two or more
532 of the shaded clades in Figure 1. In future, it might be illuminating to carry out a dat-
533 ing analysis of SE gene sequences to see if the results can help us to understand how
534 these structures might have evolved. For example, the constituents of the *sec - sel* and
535 *se27 - se28* clusters derive from the same two clades, and further these two clades are
536 two of the three clades from which the egc genes all derive. It would be interesting to

537 determine whether these common features are due to a (partially) shared inheritance
538 or whether they are coincident to independent origins.

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

545 Alongside this general picture of gene family stability and proliferation, however,
546 individual genes may also be lost. One interesting case is that of *sel29p*, observed
547 in two ST890 strains isolated in or prior to the 1940's but not seen subsequently in
548 any public sequence dataset. Could this gene have become extinct during the last 70
549 years? That it was seen only in a pseudogenised form could lend weight to such a hy-
550 pothesis, as pseudogenisation may lead to gene excision or deterioration, particularly
551 during host adaptation (31). To put the absence of copies of *sel29/sel29p* within the
552 GenBank nr/nt and RefSeq databases into sharper view, we attempted to determine
553 the Sequence Types of all 11,351 *S. aureus* genomes downloaded from RefSeq, and
554 were able to easily achieve unambiguous predictions without manual intervention in
555 11,026 cases (97%, see Data Set S2 for all predictions). We failed to find any further
556 ST890 strains within this dataset. Consequently, an alternative hypothesis of restric-
557 tion of this gene to the ST890 lineage cannot be ruled out. As well as genes that are
558 widespread across *S. aureus* strains, we have seen cases of genes restricted to a narrow
559 range of lineages (e.g. *sel33*), so this scenario would not be unprecedented. We fur-
560 ther note that the lineage itself has not become extinct, with recent reports of ST890
561 strains derived from small mammals (32, 33). However, it does not appear that any
562 of these strains have yet been subjected to whole genome sequencing. It will be inter-
563 esting to discover patterns of *sel29/sel29p* presence and absence within these strains
564 should sequencing data become available, particularly as at least one of the two NCTC
565 ST890 strains were isolated from a different host (human; see Data Set S1).

566 Our analyses indicated an association between Clonal Complex and the number
567 of SE genes/gene clusters in the strains investigated. It would be interesting to further
568 investigate possible associations between CCs and SE gene profiles (e.g. the particu-
569 lar pattern of SE genes a strain possesses), though the larger RefSeq dataset may be
570 required to achieve statistical significance. Many studies have indicated associations
571 between SE gene profiles and disease type, for example between the *egc* gene cluster
572 and both cystic fibrosis (34) and Toxic Shock Syndrome (35). Other studies have shown
573 links between disease and both SE profile and CC, such as those between CC30, infec-
574 tive endocarditis and the genes *tsst-1*, *sea*, *sed*, *see* and *sei* (36). Differences in regula-
575 tory system (7) may contribute to disease/SE gene associations, with some SE genes
576 more likely to occur in chronic rather than acute infections. The plasticity of the mobile
577 genetic elements carrying the majority of SE genes, with the different variant combi-
578 nations circulating (e.g. see Figure S2b), and the potential for their rapid loss and gain,
579 mean that selection could act swiftly on SE gene profiles. Considering such a system,
580 it would seem plausible that strains with particular SE gene profiles would be selected
581 for their roles in specific diseases and that clonal expansion would subsequently drive
582 (at least some) CCs specialised for certain diseases. Limited within-CC recombination
583 (37) would preserve these associations, establishing the patterns we see among extant
584 and preserved strains. However, SE gene co-location might also mean some observed
585 associations are indirect. A substantial meta-analysis of sequenced strains with high
586 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*, *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 colonisation/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, non-trivial 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 anti-toxins.

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 pseudogenisation 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 pseudogenisation 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 pseudogenisation vary both between genes and 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 pseudogenisation 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 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

Dataset preparation. The genome assemblies of 133 *Staphylococcus aureus* strains (see Data Set S1 for strain identities) derived from PacBio raw reads at the Wellcome Sanger Institute (WTSI) were downloaded in FASTA format from [GenBank](#). The [UniProt](#)

636 database was searched for protein and nucleotide sequences attributed to each of the
637 25 target genes (24 SEs plus *tsst-1*) from other *S. aureus* strains. In the two cases where
638 no SE matches were found (*selv* and *selw*), sequences were instead obtained from Gen-
639 Bank. Sequences were divided into two sets, where set 1 consisted of *selx* and *tsst-1*
640 and set 2 consisted of the remaining 23 SEs.

641 **Enterotoxin gene hunting.** The software tool HMMER (40) was used to build pro-
642 file hidden Markov models (pHMMs) for each gene set. The two pHMMs were then
643 used to search the 133 genome assemblies for target gene matches. Searches were
644 made using stringent parameters to guarantee full length, or close to full length, matches
645 (for set 1, $E < 1 \times e^{-10}$, $a < 5$ and $b > 600$; for set 2, $E < 1 \times e^{-10}$, $a < 88$ and $b > 590$;
646 where a is the starting co-ordinate of the match relative to the pHMM and b is the end
647 co-ordinate of the match) and more relaxed parameters ($E < 1 \times e^{-10}$ for sets 1 and 2).
648 Co-ordinates were chosen following visual inspection of the initial HMMER output for
649 a non-trivial subset of strains.

650 The HMMER accessory tool Easel was subsequently used to extract the nucleotide
651 sequences of all predicted target genes in all strains, keeping the two gene sets distinct.
652 Gene-specific pHMMs were also built for each of the 25 target genes and HMMER was
653 again used to search for and extract predicted sequences, this time separated by gene
654 identity. The gene-specific datasets were compared to the set-specific results for all
655 strains to confirm the gene family wide-approach was consistent with the gene-specific
656 approach. No inconsistencies were identified. Extracted nucleotide sequences were
657 aligned, along with 25 reference sequences (one for each target gene - see Table S1
658 for the GenBank accession numbers of all SE gene references), using MUSCLE (41) and
659 were manually divided into gene-specific groups within BioEdit (42). Using the refer-
660 ence sequences as guides, gene co-ordinates within the HMMER output were manually
661 adjusted to ensure all sequences were full length. Using the modified co-ordinates, full
662 length target gene matches were extracted from the 133 *S. aureus* strains with Easel
663 and re-aligned into gene-specific groups.

664 We also attempted both to gain support for the existence of the putative novel en-
665 terotoxin genes identified via this approach and to glean information on their origins
666 by analysing additional *S. aureus* genomes. Briefly, all 11,351 genomes within the Ref-
667 Seq database (43) available on 11/05/2020 were downloaded and were subjected to an
668 almost identical HMM-searching procedure as the 133 NCTC genome sequences. The
669 only difference in the two procedures was the use of MAFFT (44) for gene sequence
670 alignment, in place of MUSCLE, due to its ability to align tens of thousands of gene
671 sequences within a few hours (using parameters *-retree 1 -maxiterate 0 -reorder*).

672 **Phylogenetic analysis.** Translated amino acid sequences of novel or recently iden-
673 tified gene-specific groups not included in the gene hunting process, one from each
674 group, were aligned using MUSCLE along with reference sequences for previously
675 known groups and the alignment input to the IQ-TREE phylogenetic software (27) with
676 amino acid substitution model selection requested. During this process, the trans-
677 lated reference sequences for four groups (Gr5, Gr7, Gr11 and Gr12) were "repaired"
678 with minor manual editing as the group members appeared to be pseudogenes with
679 various mutations such as indel-causing frameshifts and premature stop codons. The
680 repairs, which effectively estimated the amino acid states of the sequences before
681 their putative pseudogenisation but after their divergence from the other enterotoxin
682 sequences, were made to maximise the phylogenetic signal in the dataset and hence
683 the reliability of the resulting tree. The resulting sequences are shown in Table S3.
684 We also compared the nucleotide reference sequences of all established and putative
685 gene family groups with the GenBank nr/nt database using BLAST (45).

686 **Strain typing.** Multilocus Sequence Typing (MLST) was conducted for each NCTC
687 strain using the established seven gene set for *S. aureus* (*arcC*, *aroE*, *glpF*, *gmk*, *pta*, *tpi*
688 and *yqiL*) (46). For each gene, all *S. aureus* sequences in the relevant BIGSdb database
689 (47) were downloaded and a pHMM calculated using HMMER. For each strain the MLST
690 gene sequences were extracted and concatenated into a single file, with the file sub-
691 sequently input to BIGSdb for MLST characterisation and identification of Clonal Com-
692 plex. We also carried out spa-typing for each strain, whereby the combination of differ-
693 ing repeat sequence types within the *SpA* gene was established. This process was car-
694 ried out using the `get_spa_type.py` software, which compares repeats found between
695 pairs of primer sequences against the Ridom and eGenomics typing nomenclature.
696 For 12 strains, manual editing of their genome sequence was required to achieve a
697 spa type, as sequence mutations within their *SpA* genes meant that 100% matches to
698 primer sequences were no longer achievable, preventing the software from extracting
699 the repeats. All Sequence Type, Clonal Complex and Spa Type predictions are shown
700 in Data Set S1. The Sequence Type/Clonal Complex designation process was repeated
701 for the 11,351 *S. aureus* genome assemblies downloaded from the RefSeq database.
702 Predictions for the 11,026 strains (97%) for whom manual intervention was not re-
703 quired to achieve a result are shown in Data Set S2.

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

710 **Statistical analysis.** The gene contents of strains for whom a Clonal Complex ori-
711 gin was determined were analysed along with year of isolation (set to the most recent
712 year possible given the strain metadata shown in Data Set 1) within the R statistical
713 environment (version 3.6.1) (48). A generalised linear model with a logarithmic link
714 function and Poisson error distribution was fitted to the remaining data, with the num-
715 ber of genes/gene clusters (gene clusters were counted as if they were a single gene
716 to account for the dependence of gene number counts on gene cluster presence and
717 absence) as the independent variable and Clonal Complex (a factor with six levels) and
718 year (after 1924) as dependent variables.

719 **Data availability.** All NCTC genomes have been deposited in the NCBI Sequence
720 Read Archive under NCBI BioProject accession number PRJEB6403. For individual ac-
721 cession numbers, please see Data Set S1, worksheet 1, in the supplemental material.

722 ACKNOWLEDGMENTS

723 The authors are indebted to the work of all past NCTC staff and collaborators for their
724 efforts in establishing, curating and maintaining this collection, which at time of writ-
725 ing is celebrating its 101st year of operation. We also express our gratitude to all who
726 deposited strains into the NCTC and which were examined in this manuscript. We
727 would also like to thank three anonymous reviewers for their supportive and insight-
728 ful comments. The PacBio sequencing and genome assembly of the 133 *S. aureus*
729 strains analysed within this article were funded as part of the Wellcome Trust Grant
730 no: 101503/Z/13/Z "Creation of an e-resource centre to underpin the provision and
731 use of Type and reference strains of human pathogens".

732 REFERENCES

- Graham 3rd PL, Lin SX, Larson EL. 2006. A U.S. population based survey of *Staphylococcus aureus* colonization. *Ann Intern Med* 144:318–325.
- 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.
- McCormick JK, Yarwood JM, Schlievert PM. 2001. Toxic Shock Syndrome and Bacterial Superantigens: An Update. *Annu Rev Microbiol* 55:77–104.
- 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. *The J Immunol* 171 (3):1385–1392.
- 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 (4):509–516.
- 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* 1(8228):1017–1021.
- Fisher EL, Otto M, Cheung GYC. 2018. Basis of Virulence in Enterotoxin-Mediated Staphylococcal Food Poisoning. *Front Microbiol* 9:436.
- 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 (9):e1006549.
- Argudin MA, Mendoza MC, Rodicio MR. 2010. Food Poisoning and *Staphylococcus aureus* Enterotoxins. *Toxins* 2:1751–1773.
- Thomas DY, Jarraud S, Lemerrier B, Cozon G, Echasserieu 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 (8):4724–4734.
- Okumura K, Shimomura Y, Murayama SY, Yagi J, Ubukata K, Kirikae T, Miyoshi-Akiyama T. 2012. Evolutionary paths of streptococcal and staphylococcal superantigens. *BMC Genom* 13:404.
- 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 (11):e00505–18.
- Lina G, Bohach GA, Nair SP, Hiramatsu K, Jouvin-Marche E, Mariuzza R. 2004. Standard Nomenclature for the Superantigens Expressed by *Staphylococcus*. *The J Infect Dis* 189 (12):2334–2336.
- 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 (10):E389.
- IWG-SCC. 2009. Classification of Staphylococcal Cassette Chromosome *mec* (SCC*mec*): Guidelines for Reporting Novel SCC*mec* Elements. *Antimicrob Agents Chemother* 53 (12):4961–4967.
- 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.
- 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* (SCC*mec*) Insertion Site. *The J Biol Chem* 288 (1):132–140.
- Zhang DF, Yang XY, Zhang J, Qin X, Huang X, Cui Y, Zhou M, Shi C, French NP, Shia X. 2018. Identification and characterization of two novel superantigens among *Staphylococcus aureus* complex. *Int J Med Microbiol* 308 (4):438–446.
- Noto MJ, Kreiswirth BN, Monk AB, Archer GL. 2008. Gene acquisition at the insertion site for SCC*mec*, the genomic island conferring methicillin resistance in *Staphylococcus aureus*. *J Bacteriol* 190 (4):1276–1283.
- 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 (1):65–74.
- 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.
- Jarraud S, Peyrat MA, Lim A, Tristan A, Bes M, Mougel 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*. *The J Immunol* 166 (1):669–677.
- 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 (1):38–43.
- 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 (21):6839–6849.
- 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.
- 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.
- 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 (1):268–274.
- Pushker R, Mira A, Rodríguez-Valera F. 2004. Comparative genomics of gene-family size in closely related bacteria. *Genome Biol* 5:R27.
- Sanchez-Herrero JF, Bernabeu M, Prieto A, Hüttener M, Juárez A. 2020. Gene Duplications in the Genomes of Staphylococci and Enterococci. *Front molecular biosciences* 7:160.
- 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 Reports* 5:9784.
- Goodhead I, Darby AC. 2015. Taking the pseudo out of pseudogenes. *Curr Opin Microbiol* 23:102—109.
- Monecke S, Gavier-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 (12):e0168433.
- 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.
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.
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 (8):2446–2449.
36. **Nienaber JJ, Sharma Kuinkel BK, Clarke-Pearson M, Lamlerthton 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 Jr VG, on Endocarditis-Microbiology Investigators IC.** 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 (5):704–713.
37. **Driebe EM, Sahl JW, Roe C, Bowers JR, Schupp JM, Gillece 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 (7):1–19.
38. **Aman MJ.** 2017. Superantigens of a superbug: Major culprits of *Staphylococcus aureus* disease? *Virulence* 8 (6):607–610.
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 (6):741–750.
40. **Wheeler TJ, Eddy SR.** 2013. nhmmer: DNA Homology Search With Profile HMMs. *Bioinformatics* 29:2487–2489.
41. **Edgar RC.** 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* 32 (5):1792–1797.
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 NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-White B, Ako-Adjei D, Astashyn A, Badretin A, Bao Y, Blinkova O, Brover V, Chetvernin 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, Matheron 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.** 2015. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res.* 44 (D1):D733–D745.
44. **Katoh K, Standley DM.** 2013. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 30 (4):772–780.
45. **Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ.** 1990. Basic local alignment search tool. *J. Mol. Biol.* 215:403–410.
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 (3):1008–1015.
47. **Jolley KA, Maiden MCJ.** 2010. BIGSdb: Scalable analysis of bacterial genome variation at the population level. *BMC Bioinform* 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.** 2006. Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation. *Bioinformatics* 23 (1):127–128.
50. **Bohach GA, Schlievert PM.** 1989. Conservation of the Biologically Active Portions of Staphylococcal Enterotoxins C1 and C2. *Infect Immun* 57 (7):2249–2252.

733 TABLES

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
OWWNG	0
OMI33G	0
Total	59

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

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

734 FIGURES

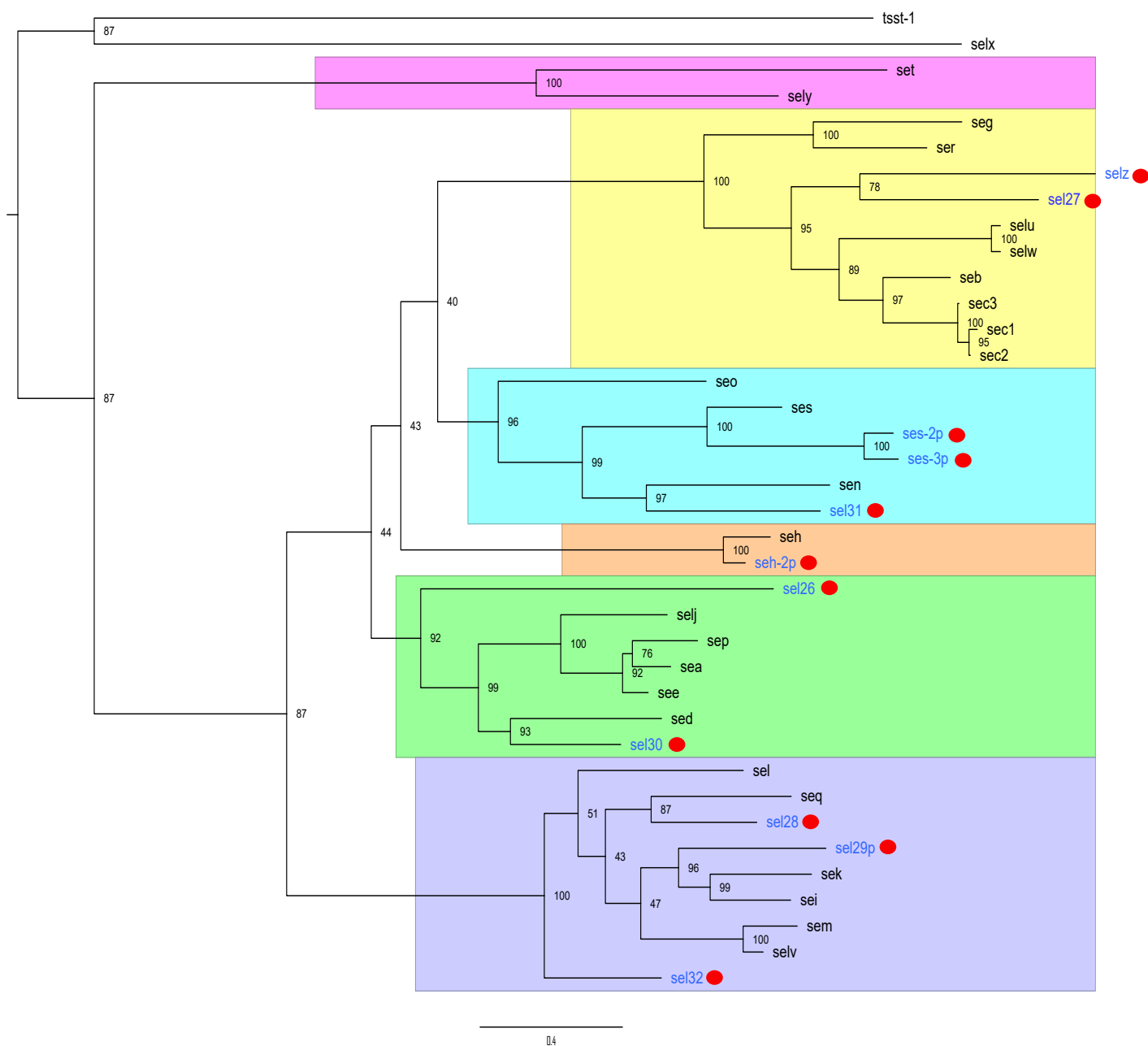


FIG 1 Phylogeny of the Staphylococcal enterotoxin genes.

Maximum likelihood (ML) phylogenetic tree of eleven of the twelve 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 coloured blocks. The tree was estimated with IQ-TREE (27) using the VT+F+R4 amino acid substitution model, maximum log-likelihood = -14507.9726, and with the *tsst-1/selx* clade used as an outgroup. 1000 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.

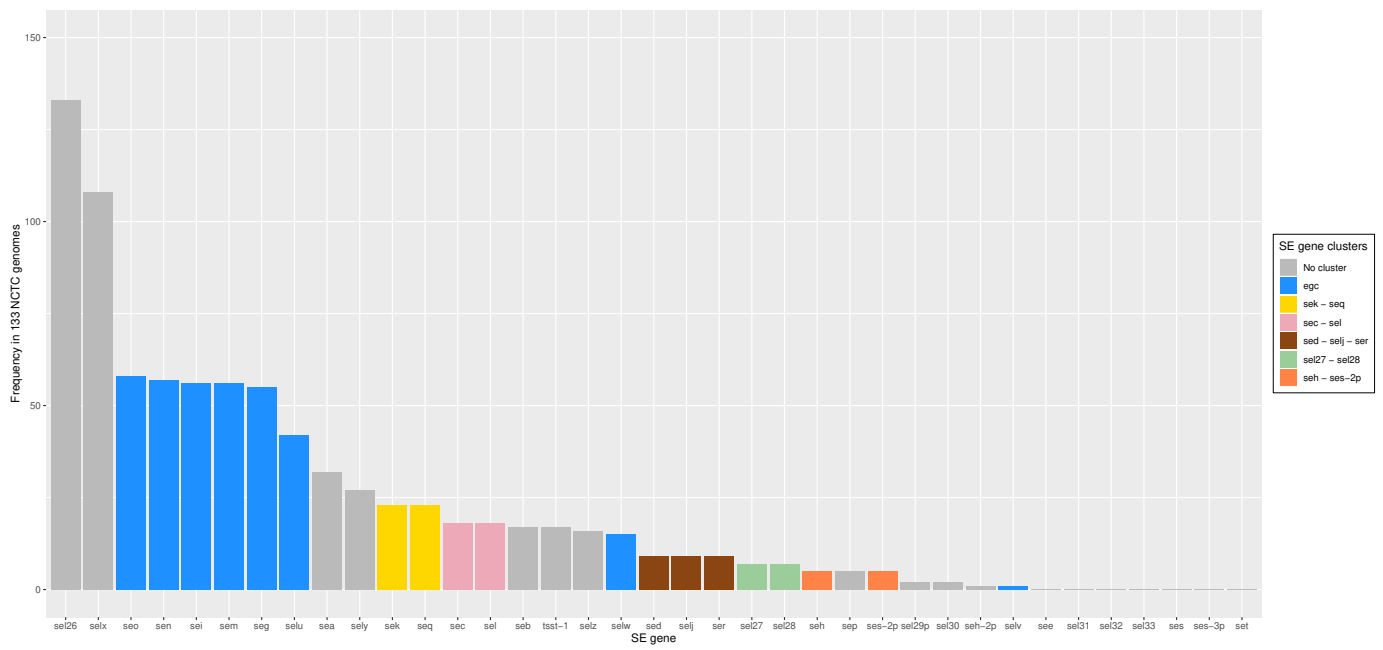


FIG 2 Staphylococcal enterotoxin gene presence. Frequencies of 36 Staphylococcal enterotoxin genes (or putative genes/pseudogenes) plus *tsst-1* in 133 NCTC *S. aureus* strains. Membership of one of the six gene clusters present in this dataset is indicated by a colour code.

Staphylococcal enterotoxin genes in NCTC genomes

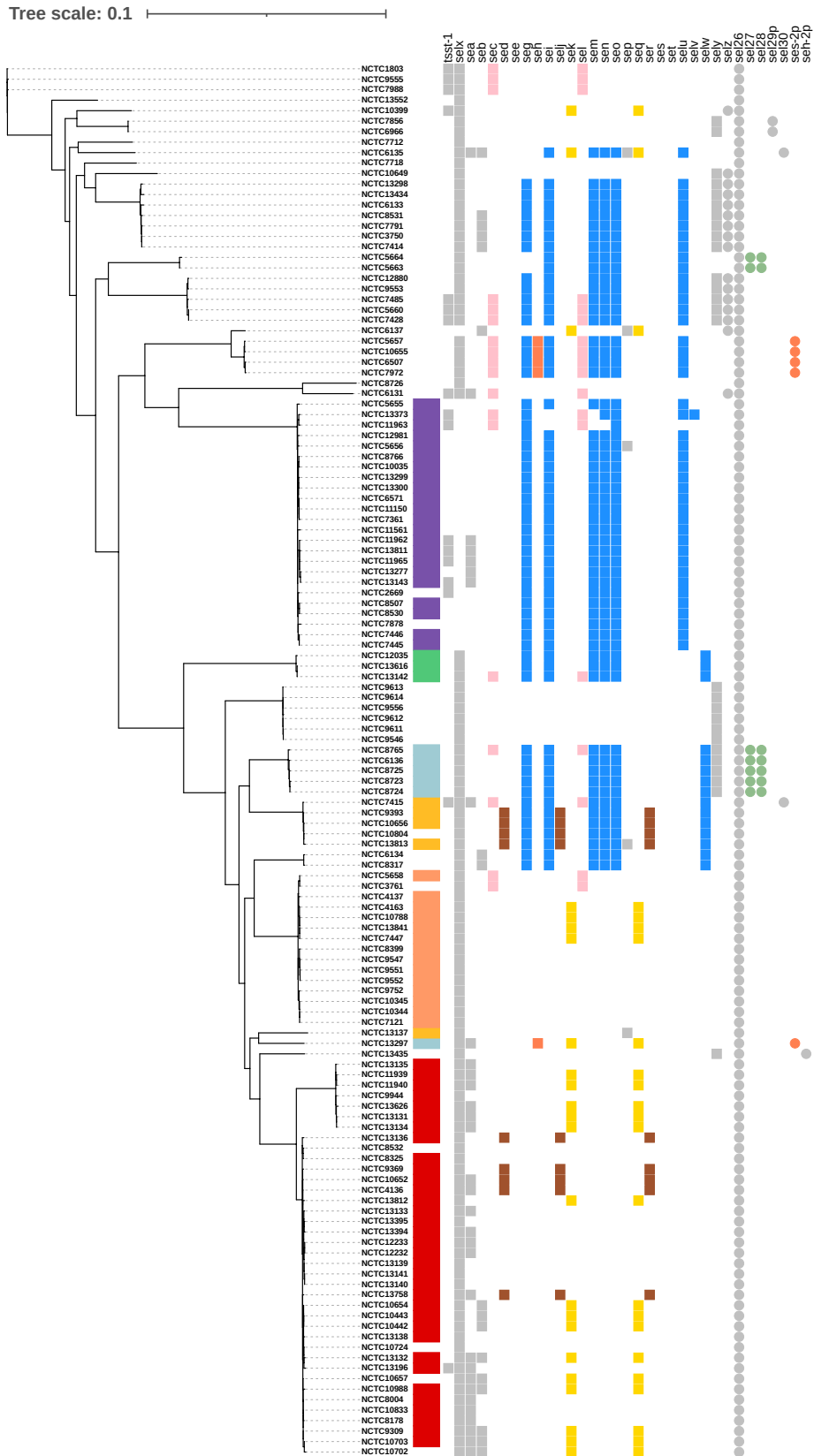


FIG 3 Staphylococcal enterotoxin gene content of 133 NCTC *Staphylococcus aureus* strains. An unrooted maximum likelihood phylogenetic tree based on 96,541 core genome SNPs is annotated with the Clonal Complex (colour strip adjacent to strain names: CC1 blue; CC5 gold; CC8 red; CC22 green; CC30 purple; CC97 orange) and SE/tsst-1 gene content (established genes as squares and recent/novel genes as circles). Gene presence is coloured 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 = -1022309.6472. The figure was generated using the [ITOL web server \(49\)](#).