1

**Average Nucleotide Identity based *Staphylococcus aureus* strain grouping allows identification of strain-specific genes in the pangenome**

3

Vishnu Raghuram[1*], Robert A Petit III[2**], Zach Karol[3], Rohan Mehta[3***], Daniel B. Weissman[3], Timothy D. Read[2****]

[1] Microbiology and Molecular Genetics Program, Graduate Division of Biological and Biomedical Sciences, Laney Graduate School, Emory University, Atlanta, Georgia, USA

[2] Division of Infectious Diseases, Department of Medicine, Emory University, Atlanta, Georgia, USA

[3] Department of Physics, Emory University, Atlanta, Georgia, USA


*Current address: Department of Clinical Microbiology, Umeå University, Umeå Sweden

**Current address: Wyoming Public Health Laboratory, Cheyenne, Wyoming USA

***Current address: Department of Biology, Elmhurst University, Elmhurst, IL, USA


****Corresponding author, Email address: tread@emory.edu


Email; ORCID
TDR: tread@emory.edu; 0000-0001-8966-9680
RAP III: robert.petit@wyo.gov; 0000-0002-1350-9426
VR: vishnu.raghuram@emory.edu; 0000-0002-7435-6435
ZK: zach.karol@emory.edu
RM: rohan.sushrut.mehta@emory.edu
DBW: daniel.weissman@emory.edu; 0000-0002-7799-1573

2

## Abstract

*Staphylococcus aureus* causes both hospital and community acquired infections in humans worldwide. Due to the high incidence of infection *S. aureus* is also one of the most sampled and sequenced pathogens today, providing an outstanding resource to understand variation at the bacterial subspecies level. We processed and downsampled 83,383 public *S. aureus* Illumina whole genome shotgun sequences and 1,263 complete genomes to produce 7,954 representative substrains. Pairwise comparison of core gene Average Nucleotide Identity (ANI) revealed a natural boundary of 99.5% that could be used to define 145 distinct strains within the species. We found that intermediate frequency genes in the pangenome (present in 10-95% of genomes) could be divided into those closely linked to strain background ("strain-concentrated") and those highly variable within strains ("strain-diffuse"). Non-core genes had different patterns of chromosome location; notably, strain-diffuse associated with prophages, strain-concentrated with the vSaβ genome island and rare genes (<10% frequency) concentrated near the origin of replication. Antibiotic genes were enriched in the strain-diffuse class, while virulence genes were distributed between strain-diffuse, strain-concentrated, core and rare classes. This study shows how different patterns of gene movement help create strains as distinct subspecies entities and provide insight into the diverse histories of important *S. aureus* functions.

## Importance

We analyzed the genomic diversity of *Staphylococcus aureus*, a globally prevalent bacterial species that causes serious infections in humans. Our goal was to build a genetic picture of the different strains of *S. aureus* and which genes may be associated with them. We used a large public dataset (>84,000 genomes) that was re-processed and subsampled to remove redundancy. We found that individual genomes could be grouped into strains by sharing > 99.5% identical nucleotide sequence of the core part of their genome. We also showed that a portion of genes that are present in intermediate frequency in the species are strongly associated with some strains but completely absent from others, suggesting a role in strain-specificity. This work lays the foundation for understanding individual gene histories of the *S. aureus* species and also outlines strategies for processing large bacterial genomic datasets.

3

## Introduction

57 *S. aureus* is a ubiquitous human pathogen capable of causing numerous disease
59 manifestations, including more than 100,000 bloodstream infections in 2017 in the US
60 alone[1]. *S. aureus* genomes typically have a ~2.8 Mbase chromosome and zero to a few
61 plasmids. Like other bacterial pathogens, its success at responding to pathogenic niches
62 comes from both adaptations in the "core" portion of the genome and non-core genes that
63 form the extended species genome, or "pangenome" [2]. Non-core genes form part of the
64 extensive genetic repertoire for evading the immune response and damaging the host
65 and have allowed *S. aureus* to survive treatment with various antibiotics developed since
66 the middle of the twentieth century [3–6].

67

68 Microbiologists have long known that there are consistent differences in phenotypes
69 between taxonomic groups below the species level in *S. aureus*. Different "strains" have
70 been shown to be more likely to cause specific disease etiologies than others. Examples
71 are Multi-Locus Sequence Type (MLST) ST582, which is associated with scalded skin
72 syndrome [7] and livestock associated CC97 infections [8]. Among other phenotypes, strains
73 also show different propensity to acquire drug resistance genes, high or low levels of
74 toxin production, and can produce different spectra of mutations when under strong
75 selection [9–12]. Understanding the genetic basis of strain-specificity therefore offers
76 potential insight into many mechanisms that define *S. aureus* pathology. Interest in strain-
77 specificity has also been prompted by attempts to use shotgun metagenomic data to
78 define environmental conditions that separate different genotypes with species [13,14].
79 However, the cardinal problem with these approaches is that there is no generally
80 accepted bacterial strain definition appropriate for the genomic era. Instead, the term
81 "strain" has been used loosely to apply to different levels of sub-species variation.

82

83 The aims of this work were to seek a consistent definition of a *S. aureus* strain that could
84 be applied to genomic and ultimately metagenomic data, to understand which portions of
85 the non-core genome were strain-associated and to survey the extent of strain variation
86 in the public data. We used an approach based on an earlier workflow [12], where we
87 reprocessed all extant public Illumina whole genome shotgun (WGS) data. Here, we
88 refined the strategy by implementing stringent steps to filter WGS potentially
89 contaminated with other bacterial contigs and *S. aureus* mixtures. We also included high-
90 quality complete genomes and dereplicated the final data set to remove very highly
91 similar sequences. Critically, we opted to define relationships between genomes based
92 on average nucleotide identity (ANI) , rather than relying on the traditional clonal complex
93 and sequence type designations of multi-locus sequence typing.

4

94 Results

95 **ANI threshold of 99.5% defines 145 *S. aureus* strains from a large public genome**
96 **dataset**

97 To get a global view of *S. aureus* genetic diversity, we used all complete genomes
98 without undefined ("N") base calls and all Illumina whole genome data sets of the species
99 available on the NCBI website in September 2022. The 83,383 whole genome data sets
100 were filtered down to 58,034 (56,771 short read genomes + 1,263 complete genomes)
101 based on having high sequence depth and quality, having no non-*S. aureus* genome
102 content, and not being potential intraspecies mixtures based on minor-allele frequency
103 (**Figure 1**, **Figure S1A; Methods**). To remove redundancy, the high-quality shotgun sets
104 and 1,263 complete genomes were clustered based on a mash distance of 0.0005
105 (approximately 50 SNPs) [12,15,16]. A randomly chosen representative of each of these 7,954
106 "substrains" was selected for downstream analysis.

107

108 The 7,954 representative substrains came from 1706 multi-locus sequence types (STs),
109 with 386 substrains not belonging to a previously assigned ST. The uneven distribution of
110 genomes across substrains and STs reflected the sampling skew towards well-known *S.*
111 *aureus* strains from predominantly clinical settings. We found that the fifteen substrains
112 that represented the most collapsed genomes, comprised 50% of the shotgun datasets.
113 The most numerous substrain, from CC22, comprised 7688 of the 58,034 whole
114 genomes (13%), while there were 5597 substrains represented by only one genome.
115 3857 out of 7,954 substrains(48%) were in ten most abundant STs (ST5, ST8, ST30,
116 ST398, ST45, ST1, ST22, ST15, ST59 and ST239), representing 39,366 out of 56,771
117 genomes (69%).

118

119 The 7,954 representative substrains were used to create a species pangenome (the
120 '7954-set'), using the PIRATE software[17] based on a minimum 50% protein sequence
121 identity. 9,533 distinct orthologous gene families were identified (we use the shortened
122 "genes" to refer to these gene families in this manuscript). Of these genes 2,008 (21.1%)
123 were considered core (found in > 95% of the genomes), 71.3% (6,794 ) were rare (<10%
124 of genomes) and 7.7% (731) were intermediate between core and rare. 90% of genes
125 were in single copy (**Figure S2**).

126

127 When pairwise average nucleotide identity (ANI) between substrains based on the
128 concatenated nucleotide sequences of the core genes (2,101,692 nt) was plotted as a
129 histogram there was a clear pattern of three strong peaks separated by distinct valleys
130 (**Figure 2A**). The left peak (smallest AN distances), we interpreted as intra-strain
131 distance, the second and third as between-strain distances within the two major *S.*
132 *aureus* clades [18], and between the clades, respectively. The threshold for intra-strain
133 relatedness appeared to be at, or very near to, 99.5%: identical to a value suggested by
134 Rodriguez-R et al to separate strains across 330 bacterial species[19]. When we used
135 99.5% as a threshold for clustering we obtained 145 groups of genomes that we termed
136 "strains" and marked each with a suffix "S99.5_" . All strain clusters had median within-
137 cluster ANI > 99.7 (**Figure S1B**). Both gene discovery rate and lineage discovery rate

5

138  were improved by dereplicating the initial 58,034 genomes compared to using a random
139  set (**Figure S1C, S1D**).

140

141  Currently, ten clonal complexes (CCs) of closely related STs are defined by the *S. aureus*
142  PubMLST site[20]. Of these, CC1, CC5, CC8, CC15, CC45, CC97 and CC121 were split
143  into 14, 3, 6, 2, 3, 5 and 5 strains, respectively, at the 99.5% clustering threshold (**Figure
144  2B**). In the case of CC1, ten strains had 7 or fewer substrains (**Figure 2C**). Two strains,
145  S99.5_9 and S99.5_36, contained substrains that had been assigned to different CCs.
146  S99.5_36 had substrains assigned CC1, CC8 and CC97 (56, 3 and 1, respectively) and
147  S99.5_9 had substrains from CC1 and CC97 (17 and 1, respectively). Substrains from
148  different CCs assigned to both S99.5_9 and S99.5_36 had at least 5 alleles in common,
149  suggesting that they were close to the threshold of being in the same CC by the rules of
150  MLST assignment (which require 5/7 common alleles). Across all strains we found that
151  >99.9% of genomes in the same strain had the same *agrD* specificity allele (1-4) of the
152  *agr* quorum sensing system (**Figure 2D**). (The one exception was strain
153  PS/BAC/317/16/W (GCF_018093225.1)[21], the single *agr* group 2 genome in 4,469 CC30
154  genomes). This result confirmed an earlier genome-based screen[15] showing that *agr* type
155  is strongly strain specific in *S. aureus*.

156

157  We noted that there was a "bump" of pairwise distances (~99.5-99.1% ANI) in the
158  otherwise clear gap between within-strain and between-strain comparisons (**Figure 2A**).
159  When we clustered substrains at 99.1% core genome ANI we found that 30 99.5%-
160  defined strains merged together to form 115 putative strains. One of the merged strains
161  comprised genomes of S99.5_2 and S99.5_27, both largely mapped to CC8. The
162  S99.5_27 strain consisted of ST239, which is known to have been created by
163  recombination of a large portion of a CC30 genome with a CC8 background [22,23]. The
164  other 9 sets of merged strains consisted of a small number of genomes. For two of the
165  merged strains, we had a complete genome which we used to align 10,000 bp sliding
166  windows against a genome from the same strain at 99.5% ANI and one from a different
167  strain that was merged at 99.9% ANI. These were strains S99.5_33 and S99.5_4 (both
168  mapped to CC45) S99.5_7 and S99.5_111 (CC15), each pair merged into one strain
169  using ANI 99.1% thresholds. Neither analysis revealed the clear pattern of large scale
170  genome replacement seen in ST239.

171

172  **Intermediate frequency genes in the pangenome can be divided into strain-**
173  **concentrated and strain-diffuse**

174  We wanted to know what proportion of the *S. aureus* accessory gene was strongly linked
175  to strain background, in the same manner as *agr* type. We adapted the commonly used
176  genetic statistic $F_{ST}$ (also known as fixation index) as a measure of segregation of a gene
177  between different strains [24]. $F_{ST}$ of 0 indicated a gene that displays no genetic
178  segregation, i.e it was indiscriminately found across different strains. In contrast, $F_{ST}$ of 1
179  indicated perfect genetic segregation, with the gene limited to all members of a group of
180  strains. Rare and core genes were constrained in their distribution and had uninformative
181  $F_{ST}$ scores around 0. Therefore we concentrated our analysis on intermediate gene
182  families.

6

183

184 Strikingly, the $F_{ST}$ distribution across intermediate genes showed a distinct bimodal
185 distribution (**Figure 3A**). This pattern disappeared when the strain labels were randomly
186 mixed and $F_{ST}$ recalculated (**Figure 3B**), reverting to a normal distribution, showing that it
187 was a feature of the specific population structure of *S. aureus* rather than an inherent
188 property of the data. From this result we divided intermediated genes into two groups
189 based on a $F_{ST}$ threshold of 0.75. Those genes with high $F_{ST}$ (296/731 (40%) intermediate
190 genes), which we termed "strain-concentrated" were strongly linked to strain
191 backgrounds, while those with low $F_{ST}$ ("strain-diffuse") (495/731 (60%) intermediate
192 genes) were more promiscuous with respect the strain background. These patterns were
193 illustrated using ten *S. aureus* toxins with a range of $F_{ST}$ scores: Leukocidins LukFS
194 (Panton Valentine Leukocidin) and LukED, Toxic Shock Syndrome toxin 1 (TSST),
195 superantigen-like protein SSL8, and different types of *Staphylococcal* Enterotoxins (SEA,
196 SEB, SEG, SEU) (**Figure 4**). Leukocidins comprise two proteins, the F component and
197 the S component, both acting synergistically to form pores in host-cell membranes [25].
198 TSST, SEs and SSL8 are superantigens or superantigen-like proteins (SAs), highly
199 potent toxins that can elicit severe inflammatory responses and other immunomodulatory
200 effects [26]. The leukocidin LukFS, enterotoxins SEA & SEB, and TSST, showed high levels
201 of gain and loss on the species tree typical of low-$F_{ST}$. In contrast, the enterotoxins SEG
202 and SEO, Leukocidin LukED, found together on genomic island vSaβ had high $F_{ST}$ (> 0.9)
203 and were either almost entirely present or absent in each strain background.

204

205 We also used $F_{ST}$ to test whether there was any association between the *agr* type of a
206 strain and intermediate gene distribution but found no similar pattern (**Figure S3**).

207

208 To investigate the differences between strain-concentrated and strain-diffuse genes
209 further in a *S. aureus* pangenome with more balanced sampling, we created the "740-
210 set", created by randomly sampling 20 shotgun assembled substrains from the most
211 common 37 strains. The 740-set had similar numbers of core and intermediate genes
212 (2,139 and 739, respectively) to the 7954-set but fewer rare genes (2,687), the latter
213 expected to increase with the number of genomes sampled in a species. The $F_{ST}$
214 distribution of the 740-set to the original pangenome was almost identical.

215 When we plotted the number of strains each gene was found in given the numbers of
216 genomes we saw two distinct patterns. The strain-concentrated genes were close to the
217 minimum possible number of strains for a given gene (dashed red line), while the strain-
218 diffuse genes were more similar to the shape of a random assortment of strains
219 (asymptotic exponential distribution; dashed blue line)(**Figure 5A**). Strain-diffuse genes
220 were present in markedly more strains at a given prevalence than strain-concentrated
221 From **Figure 5A** it was clear that rare gene distributions were extensions of the trends
222 seen in intermediate genes. These trends could not be discerned in the 7954-set
223 because the number of substrains represented in each strain was unbalanced.

224

225 **Figure 3 and 4** depict a pattern where strain-diffuse genes appeared to undergo gain and
226 loss on the phylogenetic tree at a higher rate than strain-concentrated genes. Based on

227 the results of homoplasyFinder[27] analysis on genes arrayed on the core gene phylogeny
228 of the 740-set, we found this pattern was consistent across all intermediate genes
229 (**Figure 5B**). Strain-concentrated genes mostly had fewer than 30 minimum predicted
230 state changes on the tree and there was no trend in increase of this number with
231 prevalence. Strain-diffuse genes had a higher rate of character state change, which rose
232 with prevalence initially but fell with the most common genes, probably due to saturation
233 of available state changes.

234 Because of the relatively slower rate of gene gains and losses, the strain-concentrated
235 genes contributed more to characteristic strain-specific differences in gene content than
236 strain-diffuse genes. This could be effectively visualized using t-SNE (t-distributed
237 stochastic neighbor embedding; **Figure 6**). When strain-concentrated was used as input
238 for t-SNE, the genomes that comprised individual strains were resolved into distinct
239 spatial units (**Figure 6C**). However, there was no similar pattern when strain-diffuse was
240 used (**Figure 6B**). Rare genes produced an intermediate result, with some distinctive
241 strains and some areas of the plot with mixtures of strains (**Figure 6A**) . When all non-
242 core genes were used the strains could be readily distinguished, indicating that for the t-
243 SNE approach, the strain-specific structure of strain-concentrated and rare gene content
244 was dominant to the non-strain specific strain-diffuse genes (**Figure 6D**). We also
245 visualized the effect of the different classes of non-core gene is a way that was
246 independent of strain classification: plotting the gene content similarity (represented by
247 hamming distance) of each pair of genomes against the patristic distance on the core
248 gene phylogeny (**Figure S4**). The rare and strain-diffuse genes had greater numbers of
249 gene differences between strains very closely related to each other (Patristic distance <
250 0.005) but the rate of growth of the distance in strain-concentrated genes over larger
251 distances on the phylogeny was greater. Together these results showed that strain-
252 concentrated genes provided more information about gene content differences between
253 strains than other non-core genes. We suspected that the underlying differences between
254 the two groups of genes were due to strain-concentrated genes being primarily located
255 on the chromosome and primarily spread between strains by homologous recombination,
256 whereas strain-diffuse genes were on mobile elements such as prophages, plasmids and
257 integrative conjugative elements that would be located more frequently on non-
258 chromosomal contigs. This was supported by the rate of linkage to single copy highly
259 conserved core genes (defined as whether the gene was found to be on the same contig)
260 was much lower in strain-diffuse genes (65.5%) than strain-concentrated (86.5%). By
261 comparison, the rates for rare genes were 61.5% and randomly selected genes were
262 93.5%. We used the geNomad software and database of mobile element gene [28] to see if
263 there were different distributions in the different classes of genes in the pangenome.
264 While differences between the classes were mostly statistically significant at $p < 0.05$ in
265 pairwise Tukey's tests (**Figure S5**), the difference in mean scores were mostly quite
266 small, probably reflecting the relatively small size of the *S. aureus* training set for the
267 software compared to our large pangenome sampling. The strain-diffuse genes had the
268 most distinctive signal, having the lowest mean scores for "chromosome" and "plasmid"
269 and highest for "virus". This result corroborated the association of strain-diffuse genes
270 with prophage regions of the genome.

271 We noted that the intermediate genes had a lower median clustering threshold than the
272 rare or core genes (the PIRATE software uses iterative thresholds at increasing
273 stringency to find the final clustering threshold for a gene [17]). To ensure the patterns seen

8

274    were not an artifact of lower clustering, we ran the 740-set pangenome with a minimum
275    clustering threshold of 90% amino acid identity (which we called "740-set-90"). While the
276    more stringent clustering split several rare and intermediate gene families (the "740-set-
277    90" pangenome consisted of 4,490 rare, 982 intermediate and 2,085 core) the
278    characteristic divergence in features between strain-concentrated and strain-diffuse
279    genes did not change (**Figure S6**). We also obtained similar results when the same
280    analyses were run with the original 7,954 substrain pangenome, although the unbalanced
281    nature of the collection (some strains had thousands of genomes, many only one)
282    obscured the differences between strain-concentrated and strain-diffuse in regards the
283    relationship between strains each gene was detected in at different prevalence (**Figure
284    S6A**). The strain-concentrated genes though had many fewer predicted state changes on
285    the phylogenetic tree (**Figure S6B**).

286

287    **Different non-core gene classes cluster in specific regions of the *S. aureus***
288    **chromosome, with a strong tendency for rare genes to be near the origin of**
289    **replication**

290    We used two orthologous methods to view the distribution of non-core genes on the *S.*
291    *aureus* chromosome (**Figure 7, Figure S7**). In the first method we plotted the start
292    coordinate of genes from 337 complete chromosomes(**Figure 7A, Figure S7**). There was
293    noise in the exact coordinates of individual genes but overall this method showed discrete
294    peaks in the locations of rare, and strain-concentrated and diffuse genes. The second
295    method was to link non-core genes from all 7,954 substrains to the nearest core gene on
296    the same contig (non-core genes on contigs without core genes were excluded). The
297    gross patterns of distribution of the counts of non-core genes mapped to the core nearest
298    core gene coordinate (**Figure 7B**) were similar to that in **Figure 7A**. Differences between
299    plots in the proportion of genes within each category at each genomic bin (y-axis) were
300    probably due to a combination of the indirect measurement of gene position in the linked
301    core gene method and the fact that the 7,954 substrains were are more balanced
302    reflection of *S. aureus* diversity than the 337 complete genomes.

303

304    Strain-diffuse and strain-concentrated genes had markedly distinct distributions on the
305    chromosome and were mostly located as part of distinct clusters (**Figure 7**). This could
306    also be seen clearly in the individual chromosomes of six substrains chosen to represent
307    both MRSA and MSSA from three strains (**Figure S7**). The vSaβ genome island was a
308    notably strain-concentrated-rich gene cluster, while the vSaγ island, phiSa2 and phiSa3
309    prophage were rich in strain-diffuse. The presence of strain-diffuse gene clusters was
310    more variable between genomes than strain-concentrated clusters (**Figure S7**). Some
311    genetic elements (e.g SCCmec, type VII secretion loci, phiSa1) contained a relatively
312    high proportion of both types of intermediate genes. Three regions of the chromosome
313    relatively rich in strain-concentrated genes (at approximate coordinates 100,00-300,000,
314    1,250,000-1,500,000 and 2,500,000-2,800,000) did not correspond to known genetic
315    elements , although the first region contained several genes involved in polysaccharide
316    capsule synthesis.

317

318    The high number of rare gene genes in the 0-100,000 region (which includes the

9

319 SCCmec cassette) was an outlier compared to other chromosomal regions (p-value <
320 2.2e-16, Grubbs 1-tailed test) (**Figure 7, Figure S7**). This was the case in both MRSA
321 and MSSA strains, suggesting that this region might be a hotspot for insertion of rare
322 genes, possibly through plasmid integration, rather than being specifically linked to
323 SCCmec.

324

325 **Functional differences in strain-concentrated and strain-diffuse genes**

326 $F_{ST}$ and prevalence of intermediate gene families can provide insight into ongoing
327 evolutionary processes in the species. This is illustrated by analysis of three classes of
328 genes encoding AMR (antimicrobial resistance), phage defense and virulence
329 determinants (**Figure 8**). No AMR genes[31] were found to be in the strain-concentrated
330 group but were either rare or strain-diffuse (70 (82.4%) and 15 (17.6%), respectively) .
331 This result follows from the recent introduction of many AMR genes into *S. aureus* on
332 mobile genetic elements and their frequent gains and losses below the strain level [32]. The
333 absence of fixation within strains also suggested possible loss of mobile elements in the
334 absence of antibiotic selection. Genes associated with protection from phage infection in
335 the defense-finder database [33] were mostly low prevalence (69/80 (86.3%) were rare and
336 10/80 (9.1%) intermediate had prevalence < 0.5). The low prevalence may reflect
337 diversifying selection caused by phage countermeasures. However, unlike AMR genes,
338 the majority of intermediate genes in this class were strain-concentrated, suggesting that
339 defense from phage infection may help define *S. aureus* strains. Intermediate virulence
340 genes (mostly toxins [34,35]) in the AMRFinder+ database fell into two groups: one strain-
341 diffuse with low prevalence and the other strain-concentrated with mostly higher
342 prevalence. strain-diffuse virulence genes were mostly associated with prophages and
343 Sa-PIs, while strain-concentrated genes were associated with the vSaβ genome island.
344 This partition suggested an as-yet unexplained complexity in the hierarchy of functions
345 that make up the toxin profile of an individual substrain.

10

346 Discussion

347 In this study, we distilled a starting set of >84,000 *S. aureus* genome sequences to 145
348 strains using an ANI cutoff of 99.5%, which we found to be in a natural valley between
349 clustered isolates. This threshold, or values close to it, has been reported in other studies
350 as a bacterial subspecies boundary [19]. A large number of *S. aureus* strains were rare
351 (92/145 (63.4%) represented by 1-2 substrains). While this could represent some aspect
352 of the true distribution of strain abundances in the species, it could also be a function of
353 uneven sampling of *S. aureus* genomes. There are large ascertainment biases in
354 selection as most strains are from clinical settings in western countries. It is probable that
355 the number of strains will grow significantly in the future as we extend sampling.

356

357 There is no agreed term for the highest-level bacterial subspecies level although some
358 names such as "genomevar" have been proposed [19]. We had two reasons for choosing to
359 use "strain", which is a word frequently used in microbiology but currently has a multitude
360 of different meanings . The first is to use "strain" in a way that gives it a precise definition,
361 in this case genomes that cluster together above the natural 99.5% ANI gap. The second
362 reason is that as the word is now frequently being used in metagenomic studies [13,14,36,37],
363 and by choosing "strain" to mean the highest level of subspecies, this reduces the
364 number of reference genomes needed to represent strain diversity in a species. This also
365 increases the chances of discrimination between strains using the low coverage
366 sequence read data often found in metagenome projects. However, sub-species
367 terminology needs to be formalized through standards developed by consultation with the
368 international microbiology community.

369

370 The 145 representative genomes defined here could be used for assignment of a new
371 genome to an existing strain using fastANI or similar software. If the genome was found
372 not to have >99.5% ANI to an existing strain it would be a candidate for a new strain. This
373 simple approach for strain assignment has the advantage of not needing a core
374 phylogeny calculated that is inherent to tree-based clustering and may turn out to be
375 similarly accurate owing to the population structure of the within- and between-strain
376 differences in the species (**Figure 1**). The existing MLST clonal complexes were mostly
377 mapped with a 1:1 relationship to the strains defined, and the names, which are familiar
378 in the literature, could be used as aliases for the strains. However, in some cases
379 different genome backgrounds had been designated as part of the same CC but were
380 split into more natural strain clusters by ANI. This is not surprising, as MLST schema was
381 developed for PCR amplification and sequencing, before routine whole genome
382 sequencing was available, and the seven loci used for assignment only cover a small
383 portion of the variation in the chromosome [38,39]. MLST, though useful for rapid strain
384 typing, is outperformed by whole-genome based methods for lineage assignment [39,40].

385

386 Several pangenome studies with *S. aureus* genomes have been performed for
387 epidemiological investigations [41–46], vaccine candidate discovery [47,48], and evolutionary
388 phylogenomics [49–52]. These produced a wide range of results, from 4,250 - 21,358 gene
389 total pangenome size, with cores ranging from 890 to 2,700 genes(**Table S1**). The
390 variability is a feature of the many factors that influence pangenome estimation, which

11

391    can be classed into three main groups: sample collection, data quality and bioinformatics

392    approaches. In terms of the collection, more individual genomes of a species tend to

393    produce a larger number of gene families (in an "open" pangenome) and smaller core [53].

394    Similarly, the more genetic diversity within the species increases pangenome size. We

395    used essentially all the genome data available in the public domain by Fall 2022

396    (although we ended up excluding several thousand experiments based on quality (**Figure**

397    **1**). Therefore this study probably has the largest and most diverse input *S. aureus* set

398    used to date. By reducing genome redundancy we also mitigated some of the

399    overcounting of highly sampled clones in the public databases. Ideally, all genomes for a

400    pangenome should be high-quality and complete. However, we chose to include shotgun

401    assembled genomes, which may contain a certain percentage of missing genes due to

402    contig breaks, to maximize diversity. Using shotgun assemblies also allowed us to

403    sample multiple genomes from a larger number of strains, which was important for

404    characterizing strain-diffuse and strain-concentrated genes. By reprocessing the data

405    from raw reads, we were able to filter out lower quality data and have consistent

406    assemblies (**Figure 1**). In tests, we found that pangenomes based on our shotgun

407    assemblies produce similar metrics to those estimated used only complete genomes, as

408    evidenced by the 740-set, which was composed entirely of shotgun data. For most

409    complete genomes there is no matching raw read data available in public archives, so it

410    is not possible to know whether the sequence is based on highly redundant reads

411    coverage, as it is for our Bactopia processed genomes used here. The final group of

412    factors concerns choices about bioinformatic software, and what parameters to use. Out

413    of a wide range of open source options available we chose to use highly-cited tools Bakta

414    [54] (which uses the Prodigal [55] gene finder) for annotation and PIRATE [17] for pangenome

415    estimation. PIRATE iteratively increases the threshold to report the maximum identity that

416    clusters each gene family and therefore avoids over-splitting gene families. PIRATE also

417    identifies alleles within families without creating artificial paralog gene families. Tools that

418    split paralogs into separate gene families (e.g ROARY [56] using default parameters) will

419    also produce larger numbers of gene families and fewer core genes. The choice of

420    minimum threshold for clustering proteins or genes into orthologous families (usually

421    based on percentage identity of a pairwise alignment) is important. We realized from

422    constructing the pangenome with a minimum 50% threshold that 85% of *S. aureus* genes

423    families were clustered with at least the 90% identity. When we tested the 740-set

424    pangenome with the minimum threshold increased to 90% we found a similar number of

425    core genes (2139 at 50% minimum versus 2085 at 90% minium) but the number of non-

426    core genes increased to from 3,426 to 5,472 (90%). This was because many intermediate

427    gene families had been split at the higher threshold.  However, the different threshold did

428    not affect the key result of this study was that intermediate genes could be placed into

429    two groups based on segregation with the strains defined by ANI using the $F_{ST}$ statistic.

430    Although we did not thoroughly explore different options in this study, pangenome

431    estimation in *S. aureus* could be further optimized in future benchmarking studies based

432    on the genome data collected here.

433

434    We defined three classes of *S. aureus* non-core genes with different properties. Strain-

435    diffuse genes are maintained in the population yet have a high turnover, i.e they are

436    gained and lost frequently (e.g LukFS, TSST, SEA, SEB in **Figure 4**). These genes are

437    associated with mobile elements on the chromosome such as prophages, SaPIs and

438  SCCmec and also often found on contigs unlinked to core genes, as would be expected
439  of plasmids. These genes include niche-specific functions under high selection such as
440  antibiotic resistance and certain toxins, which are classically segregated onto genetic
441  elements that undergo frequent horizontal gene transfer in bacteria. *S. aureus* strain-
442  diffuse genes are strikingly promiscuous in their strain background. Outside intra-strain
443  comparisons, there is almost no signal of phylogenetic relatedness in strain-diffuse gene
444  composition **(Figure S4)**. This suggests high rates of horizontal transfer and, over the
445  longer term, relatively weak barriers to genetic exchange compared to the strength of
446  selection for strain-diffuse genes.

447

448  The second, previously unrecognized group of intermediate genes in *S. aureus* had a
449  high $F_{ST}$ score, indicating that they segregated closely with strain core gene background.
450  Many of the genes cluster in the *S. aureus* genome islands, particularly vSaβ. The
451  elements have been described as having complex, strain-specific genetic structure [57,58].
452  Strain-concentrated genes also include significant virulence related functions located
453  outside of previously defined genetic elements such as certain type VII secretion and
454  capsule genes. strain-concentrated genes have many fewer predicted gene gains and
455  losses than strain-diffuse genes (**Figure 5**) and a much stronger phylogenetic signal
456  (**Figure S4**). This suggests that the rate of horizontal transfer of strain-diffuse genes is
457  much higher and the probable reason is that they are on self-transmissible elements such
458  as phages, plasmids (conjugative and mobilizable). The genome islands appear to have
459  evolved from prophage or SaPIs that have acquired null mutations in their genes for site-
460  specific recombination. We propose the mechanism of horizontal transfer of strain-diffuse
461  genes is indirect: homologous recombination following introduction of DNA into the donor
462  cell. Transduction is the dominant mechanism of DNA transfer in *S. aureus* and hence
463  the genes likely rely on phages and/or SaPIs for their mobility.

464

465  Rare genes probably have properties either of strain-diffuse genes (high rates of HGT) or
466  strain-concentrated genes (lower HGT rate) (**Figure 5**) but their low abundance makes
467  calculation of $F_{ST}$ the statistic meaningless. In other species (e.g *E. coli* [59]) rare genes
468  (and in some cases intermediate genes) have been reported to be strain-specific. We
469  found that rare genes had strain-specificity levels between the two classes of
470  intermediate frequency genes. In **Figure 3** some of the rare genes present in less than
471  10% of genomes are found in a significant majority 29/37 (78%) of strains. Both rare and
472  strain-diffuse genes were frequently found to be genetically linked to core genes on the
473  chromosome. While a higher proportion of strain-diffuse genes were distributed to a
474  limited number of loci, representing common insertion points for SaPIs and prophages, it
475  was a compelling finding of this study that a much higher proportion of rare genes were
476  inserted in the region near the origins of transfer (approximate coordinates 1-100,000 in
477  **Figure 7**). This was true in both MRSA and MSSA strains, hence the SCCmec element,
478  which also integrates in this region, was not solely responsible for this pattern. This
479  region of the chromosome, which is less dense in core genes, may serve as a "plasticity
480  zone" [60,61] in *S. aureus* for capture of novel genes entering the species by HGT.

481

482  This study raises two questions about the manner in which the *S. aureus* genome

13

483   evolves and the underlying selective pressures that drive the observed patterns: 1) what
484   are the forces that create the "valley" of ANI in the range of 99.1-99.5% (**Figure 1**)? and
485   2) what are the functional implications of the partitioning of intermediate genes in strain-
486   concentrated and strain-diffuse groups? The ANI valley implies that there is a limited time
487   that strains can survive as coherent taxonomic units, as measured by accumulation of
488   neutral mutations. In a recent evolutionary reconstruction, all extant *S. aureus* clonal
489   complexes tested had inferred last common ancestors in the past 250 years, most much
490   sooner[49], suggesting frequent turnover of new strains. The reasons for these replacement
491   events could be a unique historical feature of the past 2-3 centuries, caused perhaps by
492   the development of human healthcare systems and the changing chemical environment
493   of human and animal microhabitats due to technological advances but the pattern of
494   frequent strain replacement seems common to many bacterial species[19]. Possibly, strains
495   are replaced from within by the wavelike expansion of successful clones. Something like
496   this process may be happening with the expansion of USA300 since the late 1980s,
497   gradually becoming the most common CC8 strain in the USA [62,63]. This explanation
498   implies that strains occupy distinct niches, with adaptation possibly defined by the
499   composition of their non-core genes [64,65]. Substrains would then be competing with each
500   other to occupy the strain niche. There is not strong evidence of distinct within-host
501   niches for most *S. aureus* strains but there are clear associations of strains with particular
502   animal hosts[66]. New strains can also emerge from outside by genome-scale
503   recombination events, exemplified by CC239 strains, which were formed by
504   recombination of a large segment of a CC30 chromosome into a CC8 background [22,23].
505   Judging by the relatively small size of the "99.1-99.5% bump" (**Figure 1**) these types of
506   events may be a rare but ongoing process.

507

508   The second question we highlight concerns the functional implications of the partition of
509   strain-concentrated and strain-diffuse genes. There is a bias for deletion in bacterial
510   genomes[67] that implies genes maintained over time are under enduring strong selection.
511   Conversely, the strain-diffuse gene pattern can be seen as cycles of gene gain under
512   neutral selection (i.e. driven by gene transfer alone) or short term positive selection
513   followed by rapid removal. However we do not know of any studies that address the
514   underlying reasons for the difference in strain-level versus substrain-level selection.
515   Toxins are interesting in this regard because of their importance for *S. aureus* virulence.
516   Why are some toxins maintained as core functions (e.g alpha-toxin (*hly)*), some strain-
517   concentrated (e.g enterotoxin G (*seg*)) and some strain-diffuse, present in diverse
518   substrains (e.g Panton-Valentine leukocidin (*lukFS*))? (**Figure 4**). The superantigen-type
519   toxins are split between strain-concentrated and strain-diffuse genes, suggesting that
520   former functions may be strongly linked to strain niches. Related to these issues is the
521   question of long-term maintenance of diversity of strain-concentrated genes under
522   conditions of relatively low transfer rate and rapid strain extinction that would suggest a
523   high rate of stochastic loss. Could there be frequency-dependent selection operating
524   across the *S. aureus* species on strain-concentrated genes?.

525

526   In summary, this work revealed a new partition in the structure of the *S. aureus*
527   pangenome that will spur further studies on genome evolution and subspeciation in the
528   species. The methodology for refining large amounts of public data, defining strains using

14

529　ANI and following strain-specificity of the pangenome using $F_{st}$ can also be applied to
530　other bacterial species.  Comparisons to other species, particularly from the
531　*Staphylococcus* genus, will reveal the commonalities and unique selective pressures
532　acting on the pangenome of this dangerous pathogen.

533

534　<u>Acknowledgements</u>

15

## Methods

### Public genome collection, processing and filtering

Bactopia v1.7.0 was used to download and process all genomes used in this dataset. Bactopia is a software pipeline for comprehensive analysis of bacterial genomes based on Nextflow [68,69]. The command "*bactopia search "Staphylococcus aureus" --prefix saureus*" was used to download all *S. aureus* short-read sequences available on Sequence Read Archive (SRA) as of September 2022. Bactopia used SKESA to assemble genomes, Bakta to annotate and Snippy for variant calling [70,71]. Assembly quality was evaluated using QUAST and CheckM [72,73]. *S. aureus* CC and ST were based on the pubmlst database [20]. (https://pubmlst.org/bigsdb? db=pubmlst_saureus_seqdef&page=downloadProfiles&scheme_id=1). AgrVATE v1.0.5 was used to assign *agr* types [15]. Only samples having greater than 50× coverage, mean per-read quality greater than 20, mean read length greater than 75 bp, and an assembly with less than 200 contigs were considered for the analysis (corresponding to 'Gold' and 'Silver' ranks as designated by Bactopia. Samples that were detected as not *S. aureus* according to kmer based identification or CheckM were then removed. Coverage for all samples were capped at 100x. For every sample, bactopia performs variant calling using Snippy against an auto-chosen reference sequence based on the smallest MASH distance to a complete *S. aureus* genome in RefSeq [70,74]. For each variant identified, the allele frequencies were calculated from the bam files using bcftools mpileup [75]. Samples having average minor allele frequency > 0.05 were considered mixed strains and therefore removed. Samples having total number of variants > 150,000 compared to the auto-chosen reference (or more than 5% of the genome) were also considered non-*S. aureus* and removed [76]. This process reduced 83,383 samples to 56,771. Since Bactopia collected and processed only short read *S. aureus* data, we added complete *S. aureus* genome sequences to this set. Out of 1,475 complete genomes publicly available as of February 2023, 1,263 did not have any 'N' characters in their assemblies and were added to the filtered dataset of 56,771, leading to a total of 58,034 genomes. The 212 complete genomes containing 'N' characters were not used in this study.

### Substrain dereplication

Samples were grouped by their MLST types as assigned by Bactopia and for each ST, an all vs all MASH distance estimation [74] was run. Samples with a MASH distance < 0.0005 were grouped into clusters and a random genome was chosen as the cluster representative [16]. However, where possible, we used complete genomes as the cluster representative. Samples with unassigned STs were grouped together and treated the same. The resulting final dereplicated set comprised 7954 genomes and was used for pangenome construction.

### Pangenome analysis

The bakta annotation produced by the original Bactopia run was used as input for pangenome estimation with PIRATE 1.0.5 [17]. PIRATE was run using default parameters with the additional flags -a to obtain core genome alignments and -k "--diamond" to use DIAMOND for the amino-acid sequence comparisons [77]. SNP-sites v 2.5.1 [78] was run on the PIRATE core genome alignment to extract only polymorphic sites (709,911 sites) and the resulting alignment was used to construct a core genome phylogeny with FastTree v 2.1.11 [79](GTR model, 1000 bootstrap resamples). The phylogeny was visualized using

16

588  the R package ggtree [80,81]. We used Homoplasyfinder[27] to count the number of state
589  changes of each non-core gene on the phylogeny. geNomad v1.5[28] was used to predict
590  mobile genetic elements.

**Strain definition based on ANI**

592  All-vs-All pairwise ANI was calculated for the 7,954 dereplicated genomes using fastANI
593  v1.33 [76]. Strain assignments were performed based on average linkage hierarchical
594  clustering and samples that had ANI 99.5% or greater were clustered together. The
595  average ANI of each genome with every other genome in a given cluster was calculated
596  and the genome with the highest average ANI was assigned as the strain representative.

**Calculating $F_{ST}$**

598  We created a custom R function to calculate the $F_{ST}$ for each gene, with group membership
599  defined as strain type, clonal complex or *agr* group, depending on the purpose of the
600  comparison. The input was a binary presence/absence data frame, with genes as columns
601  and genomes as rows. $F_{ST}$ was calculated using Weir's formula [24].

**Creating the 740-set and 740-set-90 pangenomes**

603  We randomly subsampled the 37 strains with > 20 substrains to 20 substrain genomes
604  each. We rerun PIRATE 1.0.5 with default parameters and created a core pangnome tree
605  using FastTree v 2.1.11 as described above. To create the "740-set-90" pangenome we
606  the 740 genomes through PIRATE 1.0.5 with minimum clustering threshold of 90% amino
607  acid identity.

**Chromosomal locations of non-core genes**

609  We used two methods for mapping chromosomal locations of non-core genes based on
610  the co-ords output of the PIRATE 1.0.5 pipeline for the 7954-set and 740-set
611  pangenomes. First we screened 377 complete substrain genome that had *dnaA* as their
612  first gene by BLAST and collated the start coordinate of each non-core gene. The second
613  method was to collate the start coordinate of nearest core gene on the same contig as
614  each non-core gene. For each class of non-core gene 20,000 random genes were
615  selected as well as a control of 20,000 genes of all classes (including core). If the non-
616  core gene was on a contig that did not have a core gene then its status was returned as
617  "unlinked".

**Antibiotic resistance, virulence and phage defense functions**

619  To assign antibiotic-resistance genes we queried representative protein sequences of
620  each gene family of the 7954-set produced by PIRATE against the AMRFinder+[31]
621  database using tblastn [82] with a threshold of >= 90% identity as a match. We filtered the
622  out virulence-associated genes using matches the terms: "serine_protease",
623  "enterotoxin", "hemolysin", "Panton", "adhesin", "complement", "aureolysin", "exfoliative",
624  "toxin", "intracellular_survival", "serum_survival" and "leukocidin" and the kept the
625  remainder as antibiotic-resistance gene matches. To assign phage defense related
626  functions, we queried the 7954-set representative proteins against the online
627  defensefinder database[33] (https://defense-finder.mdmparis-lab.com/) on 2023-10-17.

**Statistical analysis and data visualization**

629  All statistics and tSNE were performed in R using package rstatix [83]. All plots were
630  visualized using R package ggplot2 [84]. Other visualizations were performed using draw.io
631  and Sakneymatic [85,86].

17

## Data availability

PIRATE pangenome outputs, genes and strain lists and representative genome sets are available on Zenodo https://zenodo.org/records/10471309.

18

635 References

636  1   Kourtis AP, Hatfield K, Baggs J, Mu Y, See I, Epson E, *et al.* Vital Signs: Epidemiology
637      and Recent Trends in Methicillin-Resistant and in Methicillin-Susceptible
638      Staphylococcus aureus Bloodstream Infections - United States. *MMWR Morb Mortal*
639      *Wkly Rep* 2019;**68**:214–9. https://doi.org/10.15585/mmwr.mm6809e1.
640  2   Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, *et al.* Genome
641      analysis of multiple pathogenic isolates of Streptococcus agalactiae: implications for the
642      microbial 'pan-genome'. *Proc Natl Acad Sci U S A* 2005;**102**:13950–5.
643      https://doi.org/10.1073/pnas.0506758102.
644  3   Howden BP, Giulieri SG, Wong Fok Lung T, Baines SL, Sharkey LK, Lee JYH, *et al.*
645      Staphylococcus aureus host interactions and adaptation. *Nat Rev Microbiol* 2023.
646      https://doi.org/10.1038/s41579-023-00852-y.
647  4   Vestergaard M, Frees D, Ingmer H. Antibiotic Resistance and the MRSA Problem.
648      *Microbiol Spectr* 2019;**7**.: https://doi.org/10.1128/microbiolspec.GPP3-0057-2018.
649  5   Peschel A, Otto M. Phenol-soluble modulins and staphylococcal infection. *Nat Rev*
650      *Microbiol* 2013;**11**:667–73. https://doi.org/10.1038/nrmicro3110.
651  6   Spaan AN, van Strijp JAG, Torres VJ. Leukocidins: staphylococcal bi-component pore-
652      forming toxins find their receptors. *Nat Rev Microbiol* 2017.
653      https://doi.org/10.1038/nrmicro.2017.27.
654  7   Azarian T, Cella E, Baines SL, Shumaker MJ, Samel C, Jubair M, *et al.* Genomic
655      Epidemiology and Global Population Structure of Exfoliative Toxin A-Producing
656      Staphylococcus aureus Strains Associated With Staphylococcal Scalded Skin
657      Syndrome. *Front Microbiol* 2021;**12**:2307. https://doi.org/10.3389/fmicb.2021.663831.
658  8   Spoor LE, McAdam PR, Weinert LA, Rambaut A, Hasman H, Aarestrup FM, *et al.*
659      Livestock origin for a human pandemic clone of community-associated methicillin-
660      resistant Staphylococcus aureus. *MBio* 2013;**4**.: https://doi.org/10.1128/mBio.00356-13.
661  9   Su M, Lyles JT, Petit RA Iii, Peterson J, Hargita M, Tang H, *et al.* Genomic analysis of
662      variability in Delta-toxin levels between Staphylococcus aureus strains. *PeerJ*
663      2020;**8**:e8717. https://doi.org/10.7717/peerj.8717.
664 10   Benson MA, Ohneck EA, Ryan C, Alonzo F 3rd, Smith H, Narechania A, *et al.* Evolution
665      of hypervirulence by a MRSA clone through acquisition of a transposable element. *Mol*
666      *Microbiol* 2014;**93**:664–81. https://doi.org/10.1111/mmi.12682.
667 11   Su M, Davis MH, Peterson J, Solis-Lemus C, Satola SW, Read TD. Effect of genetic
668      background on the evolution of Vancomycin-Intermediate Staphylococcus aureus
669      (VISA). *PeerJ* 2021;**9**:e11764. https://doi.org/10.7717/peerj.11764.
670 12   Petit RA 3rd, Read TD. Staphylococcus aureus viewed from the perspective of 40,000+
671      genomes. *PeerJ* 2018;**6**:e5261. https://doi.org/10.7717/peerj.5261.
672 13   Van Rossum T, Ferretti P, Maistrenko OM, Bork P. Diversity within species: interpreting
673      strains in microbiomes. *Nat Rev Microbiol* 2020. https://doi.org/10.1038/s41579-020-
674      0368-1.
675 14   Liao H, Ji Y, Sun Y. High-resolution strain-level microbiome composition analysis from
676      short reads. *Microbiome* 2023;**11**:183. https://doi.org/10.1186/s40168-023-01615-w.
677 15   Raghuram V, Alexander AM, Loo HQ, Petit RA 3rd, Goldberg JB, Read TD. Species-
678      Wide Phylogenomics of the Staphylococcus aureus Agr Operon Revealed Convergent
679      Evolution of Frameshift Mutations. *Microbiol Spectr* 2022:e0133421.
680      https://doi.org/10.1128/spectrum.01334-21.
681 16   Raghuram V, Read T. *Help, I have too many genome sequences!*. 2022.
682      https://doi.org/10.5281/zenodo.7278310.
683 17   Bayliss SC, Thorpe HA, Coyle NM, Sheppard SK, Feil EJ. PIRATE: A fast and scalable
684      pangenomics toolbox for clustering diverged orthologues in bacteria. *Gigascience*
685      2019;**8**:598391. https://doi.org/10.1093/gigascience/giz119.
686 18   Planet PJ, Narechania A, Chen L, Mathema B, Boundy S, Archer G, *et al.* Architecture
687      of a Species: Phylogenomics of Staphylococcus aureus. *Trends Microbiol* 2016.
688      https://doi.org/10.1016/j.tim.2016.09.009.

19

689 19 Rodriguez-R LM, Conrad RE, Viver T, Feistel DJ, Lindner BG, Venter SN, *et al.* An ANI
690     gap within bacterial species that advances the definitions of intra-species units. *MBio*
691     2023:e0269623. https://doi.org/10.1128/mbio.02696-23.
692 20 Jolley KA, Bray JE, Maiden MCJ. Open-access bacterial population genomics: BIGSdb
693     software, the PubMLST.org website and their applications. *Wellcome Open Res*
694     2018;**3**:124. https://doi.org/10.12688/wellcomeopenres.14826.1.
695 21 Yousuf B, Flint A, Weedmark K, McDonald C, Bearne J, Pagotto F, *et al.* Genome
696     Sequence of Staphylococcus aureus Strain PS/BAC/317/16/W, Isolated from
697     Contaminated Platelet Concentrates in England. *Microbiol Resour Announc*
698     2021;**10**:e0057721. https://doi.org/10.1128/MRA.00577-21.
699 22 Gill JL, Hedge J, Wilson DJ, MacLean RC. Evolutionary Processes Driving the Rise and
700     Fall of Staphylococcus aureus ST239, a Dominant Hybrid Pathogen. *MBio*
701     2021;**12**:e0216821. https://doi.org/10.1128/mBio.02168-21.
702 23 Robinson DA, Enright MC. Evolution of Staphylococcus aureus by large chromosomal
703     replacements. *J Bacteriol* 2004;**186**:1060–4.
704 24 Weir BS. Estimating F-statistics: A historical view. *Philos Sci* 2012;**79**:637–43.
705     https://doi.org/10.1086/667904.
706 25 Melles DC, van Leeuwen WB, Boelens HAM, Peeters JK, Verbrugh HA, van Belkum A.
707     Panton-Valentine leukocidin genes in Staphylococcus aureus. *Emerg Infect Dis*
708     2006;**12**:1174–5. https://doi.org/10.3201/eid1207.050865.
709 26 Krakauer T. Staphylococcal Superantigens: Pyrogenic Toxins Induce Toxic Shock.
710     *Toxins* 2019;**11.**: https://doi.org/10.3390/toxins11030178.
711 27 Crispell J, Balaz D, Gordon SV. HomoplasyFinder: a simple tool to identify homoplasies
712     on a phylogeny. *Microb Genom* 2019;**5.**: https://doi.org/10.1099/mgen.0.000245.
713 28 Camargo AP, Roux S, Schulz F, Babinski M, Xu Y, Hu B, *et al.* Identification of mobile
714     genetic elements with geNomad. *Nat Biotechnol* 2023. https://doi.org/10.1038/s41587-
715     023-01953-y.
716 29 Gill SR, Fouts DE, Archer GL, Mongodin EF, Deboy RT, Ravel J, *et al.* Insights on
717     evolution of virulence and resistance from the complete genome analysis of an early
718     methicillin-resistant Staphylococcus aureus strain and a biofilm-producing methicillin-
719     resistant Staphylococcus epidermidis strain. *J Bacteriol* 2005;**187**:2426–38.
720     https://doi.org/10.1128/JB.187.7.2426-2438.2005.
721 30 Warne B, Harkins CP, Harris SR, Vatsiou A, Stanley-Wall N, Parkhill J, *et al.* The
722     Ess/Type VII secretion system of Staphylococcus aureus shows unexpected genetic
723     diversity. *BMC Genomics* 2016;**17**:222. https://doi.org/10.1186/s12864-016-2426-7.
724 31 Feldgarden M, Brover V, Gonzalez-Escalona N, Frye JG, Haendiges J, Haft DH, *et al.*
725     AMRFinderPlus and the Reference Gene Catalog facilitate examination of the genomic
726     links among antimicrobial resistance, stress response, and virulence. *Sci Rep*
727     2021;**11**:12728. https://doi.org/10.1038/s41598-021-91456-0.
728 32 Chambers HF, Deleo FR. Waves of resistance: Staphylococcus aureus in the antibiotic
729     era. *Nat Rev Microbiol* 2009;**7**:629–41. https://doi.org/10.1038/nrmicro2200.
730 33 Tesson F, Hervé A, Mordret E, Touchon M, d'Humières C, Cury J, *et al.* Systematic and
731     quantitative view of the antiviral arsenal of prokaryotes. *Nat Commun* 2022;**13**:2561.
732     https://doi.org/10.1038/s41467-022-30269-9.
733 34 Xia G, Wolz C. Phages of Staphylococcus aureus and their impact on host evolution.
734     *Infect Genet Evol* 2014;**21**:593–601. https://doi.org/10.1016/j.meegid.2013.04.022.
735 35 McCarthy AJ, Lindsay JA. The distribution of plasmids that carry virulence and
736     resistance genes in Staphylococcus aureus is lineage associated. *BMC Microbiol*
737     2012;**12**:104. https://doi.org/10.1186/1471-2180-12-104.
738 36 Beghini F, McIver LJ, Blanco-Miguez A, Dubois L, Asnicar F, Maharjan S, *et al.*
739     Integrating taxonomic, functional, and strain-level profiling of diverse microbial
740     communities with bioBakery 3. *Cold Spring Harbor Laboratory*
741     2020:2020.11.19.388223. https://doi.org/10.1101/2020.11.19.388223.
742 37 Jin X, Yu FB, Yan J, Weakley AM, Dubinkina V, Meng X, *et al.* Culturing of a complex
743     gut microbial community in mucin-hydrogel carriers reveals strain- and gene-associated

20

744   spatial organization. *Nat Commun* 2023;**14**:3510. https://doi.org/10.1038/s41467-023-
745   39121-0.
746   38   Maiden MC, Bygraves JA, Feil E, Morelli G, Russell JE, Urwin R, *et al.* Multilocus
747   sequence typing: a portable approach to the identification of clones within populations of
748   pathogenic microorganisms. *Proc Natl Acad Sci U S A* 1998;**95**:3140–5.
749   https://doi.org/10.1073/pnas.95.6.3140.
750   39   Maiden MCJ, van Rensburg MJJ, Bray JE, Earle SG, Ford SA, Jolley KA, *et al.* MLST
751   revisited: the gene-by-gene approach to bacterial genomics. *Nat Rev Microbiol*
752   2013;**11**:728–36. https://doi.org/10.1038/nrmicro3093.
753   40   Falush D. Toward the use of genomics to study microevolutionary change in bacteria.
754   *PLoS Genet* 2009;**5**:e1000627. https://doi.org/10.1371/journal.pgen.1000627.
755   41   Jamrozy DM, Harris SR, Mohamed N, Peacock SJ, Tan CY, Parkhill J, *et al.* Pan-
756   genomic perspective on the evolution of the Staphylococcus aureus USA300 epidemic.
757   *Microb Genom* 2016;**2**:e000058. https://doi.org/10.1099/mgen.0.000058.
758   42   Long DR, Wolter DJ, Lee M, Precit M, McLean K, Holmes E, *et al.* Polyclonality, Shared
759   Strains, and Convergent Evolution in Chronic CF S. aureus Airway Infection. *Am J*
760   *Respir Crit Care Med* 2020. https://doi.org/10.1164/rccm.202003-0735OC.
761   43   Montelongo C, Mores CR, Putonti C, Wolfe AJ, Abouelfetouh A. Whole-Genome
762   Sequencing of Staphylococcus aureus and Staphylococcus haemolyticus Clinical
763   Isolates from Egypt. *Microbiol Spectr* 2022;**10**:e0241321.
764   https://doi.org/10.1128/spectrum.02413-21.
765   44   Xu Z, Yuan C. Molecular Epidemiology of Staphylococcus aureus in China Reveals the
766   Key Gene Features Involved in Epidemic Transmission and Adaptive Evolution.
767   *Microbiol Spectr* 2022:e0156422. https://doi.org/10.1128/spectrum.01564-22.
768   45   Holm MKA, Jørgensen KM, Bagge K, Worning P, Pedersen M, Westh H, *et al.*
769   Estimated Roles of the Carrier and the Bacterial Strain When Methicillin-Resistant
770   Staphylococcus aureus Decolonization Fails: a Case-Control Study. *Microbiol Spectr*
771   2022:e0129622. https://doi.org/10.1128/spectrum.01296-22.
772   46   Cella E, Sutcliffe CG, Tso C, Paul E, Ritchie N, Colelay J, *et al.* Carriage prevalence
773   and genomic epidemiology of Staphylococcus aureus among Native American children
774   and adults in the Southwestern USA. *Microbial Genomics* 2022;**8**:000806.
775   https://doi.org/10.1099/mgen.0.000806.
776   47   Naz K, Ullah N, Naz A, Irum S, Dar HA, Zaheer T, *et al.* The Epidemiological and
777   Pangenome Landscape of Staphylococcus aureus and Identification of Conserved
778   Novel Candidate Vaccine Antigens. *Curr Proteomics* 2022;**19**:114–26.
779   https://doi.org/10.2174/1570164618666210212122847.
780   48   Naz K, Naz A, Ashraf ST, Rizwan M, Ahmad J, Baumbach J, *et al.* PanRV: Pangenome-
781   reverse vaccinology approach for identifications of potential vaccine candidates in
782   microbial pangenome. *BMC Bioinformatics* 2019;**20**:123.
783   https://doi.org/10.1186/s12859-019-2713-9.
784   49   Yebra G, Harling-Lee JD, Lycett S, Aarestrup FM, Larsen G, Cavaco LM, *et al.*
785   Multiclonal human origin and global expansion of an endemic bacterial pathogen of
786   livestock. *Proc Natl Acad Sci U S A* 2022;**119**:e2211217119.
787   https://doi.org/10.1073/pnas.2211217119.
788   50   Blaustein RA, McFarland AG, Ben Maamar S, Lopez A, Castro-Wallace S, Hartmann
789   EM. Pangenomic Approach To Understanding Microbial Adaptations within a Model
790   Built Environment, the International Space Station, Relative to Human Hosts and Soil.
791   *mSystems* 2019;**4**:e00281–18. https://doi.org/10.1128/mSystems.00281-18.
792   51   Rao RT, Sivakumar N, Jayakumar K. Analyses of Livestock-Associated Staphylococcus
793   aureus Pan-Genomes Suggest Virulence Is Not Primary Interest in Evolution of Its
794   Genome. *OMICS* 2019;**23**:224–36. https://doi.org/10.1089/omi.2019.0005.
795   52   John J, George S, Nori SRC, Nelson-Sathi S. Evolutionary route of resistant genes in
796   Staphylococcus aureus. *Genome Biol Evol* 2019. https://doi.org/10.1093/gbe/evz213.
797   53   Vernikos G, Medini D, Riley DR, Tettelin H. Ten years of pan-genome analyses. *Curr*
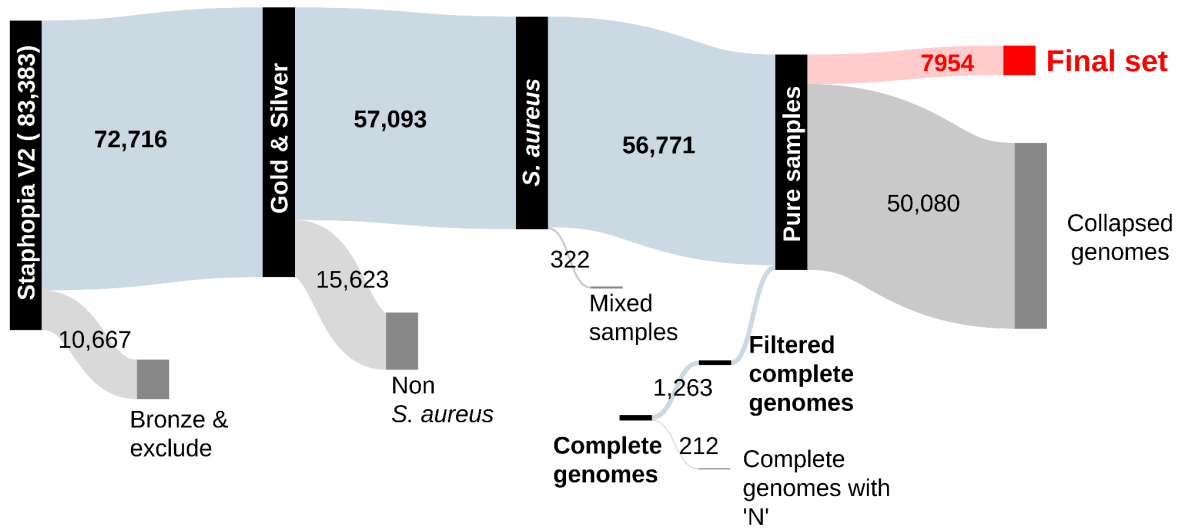798   *Opin Microbiol* 2014;**23C**:148–54. https://doi.org/10.1016/j.mib.2014.11.016.

799   54   Schwengers O, Jelonek L, Dieckmann MA, Beyvers S, Blom J, Goesmann A. Bakta:
800        rapid and standardized annotation of bacterial genomes via alignment-free sequence
801        identification. *Microb Genom* 2021;**7.**: https://doi.org/10.1099/mgen.0.000685.
802   55   Hyatt D, Chen G-L, Locascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:
803        prokaryotic gene recognition and translation initiation site identification. *BMC
804        Bioinformatics* 2010;**11**:119. https://doi.org/10.1186/1471-2105-11-119.
805   56   Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, *et al.* Roary: rapid
806        large-scale prokaryote pan genome analysis. *Bioinformatics* 2015;**31**:3691–3.
807        https://doi.org/10.1093/bioinformatics/btv421.
808   57   Baba T, Bae T, Schneewind O, Takeuchi F, Hiramatsu K. Genome sequence of
809        Staphylococcus aureus strain Newman and comparative analysis of staphylococcal
810        genomes: polymorphism and evolution of two major pathogenicity islands. *J Bacteriol*
811        2008;**190**:300–10. https://doi.org/10.1128/JB.01000-07.
812   58   Kläui AJ, Boss R, Graber HU. Characterization and Comparative Analysis of the
813        Staphylococcus aureus Genomic Island vSaβ - an in silico Approach. *J Bacteriol* 2019.
814        https://doi.org/10.1128/JB.00777-18.
815   59   Horesh G, Taylor-Brown A, McGimpsey S, Lassalle F, Corander J, Heinz E, *et al.*
816        Different evolutionary trends form the twilight zone of the bacterial pan-genome. *Microb
817        Genom* 2021;**7**:2021.02.15.431222. https://doi.org/10.1099/mgen.0.000670.
818   60   Read TD, Brunham RC, Shen C, Gill SR, Heidelberg JF, White O, *et al.* Genome
819        sequences of Chlamydia trachomatis MoPn and Chlamydia pneumoniae AR39. *Nucleic
820        Acids Res* 2000;**28**:1397–406.
821   61   Dobrindt U, Hacker J. Whole genome plasticity in pathogenic bacteria. *Curr Opin
822        Microbiol* 2001;**4**:550–7. https://doi.org/10.1016/s1369-5274(00)00250-2.
823   62   David MZ, Daum RS. Community-associated methicillin-resistant Staphylococcus
824        aureus: epidemiology and clinical consequences of an emerging epidemic. *Clin
825        Microbiol Rev* 2010;**23**:616–87. https://doi.org/10.1128/CMR.00081-09.
826   63   Diekema DJ, Richter SS, Heilmann KP, Dohrn CL, Riahi F, Tendolkar S, *et al.*
827        Continued emergence of USA300 methicillin-resistant Staphylococcus aureus in the
828        United States: results from a nationwide surveillance study. *Infect Control Hosp
829        Epidemiol* 2014;**35**:285–92. https://doi.org/10.1086/675283.
830   64   Vos M, Eyre-Walker A. Are pangenomes adaptive or not? *Nat Microbiol* 2017:1576.
831        https://doi.org/10.1038/s41564-017-0067-5.
832   65   McInerney JO, McNally A, O'Connell MJ. Why prokaryotes have pangenomes. *Nature
833        Microbiology* 2017;**2**:17040. https://doi.org/10.1038/nmicrobiol.2017.40.
834   66   Richardson EJ, Bacigalupe R, Harrison EM, Weinert LA, Lycett S, Vrieling M, *et al.*
835        Gene exchange drives the ecological success of a multi-host bacterial pathogen. *Nat
836        Ecol Evol* 2018. https://doi.org/10.1038/s41559-018-0617-0.
837   67   Mira A, Ochman H, Moran NA. Deletional bias and the evolution of bacterial genomes.
838        *Trends Genet* 2001;**17**:589–96.
839   68   Petit RA, Read TD. Bactopia: a Flexible Pipeline for Complete Analysis of Bacterial
840        Genomes. *mSystems* 2020;**5.**: https://doi.org/10.1128/mSystems.00190-20.
841   69   Di Tommaso P, Chatzou M, Floden EW, Barja PP, Palumbo E, Notredame C. Nextflow
842        enables reproducible computational workflows. *Nat Biotechnol* 2017;**35**:316–9.
843        https://doi.org/10.1038/nbt.3820.
844   70   Seemann T. *snippy: Rapid haploid variant calling and core genome alignment*. 2023.
845        URL: https://github.com/tseemann/snippy (Accessed 25 September 2023).
846   71   Souvorov A, Agarwala R, Lipman DJ. SKESA: strategic k-mer extension for scrupulous
847        assemblies. *Genome Biol* 2018;**19**:153. https://doi.org/10.1186/s13059-018-1540-z.
848   72   Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for
849        genome assemblies. *Bioinformatics* 2013;**29**:1072–5.
850        https://doi.org/10.1093/bioinformatics/btt086.
851   73   Parks DH, Imelfort M, Skennerton CT, Hugenholtz P, Tyson GW. CheckM: assessing
852        the quality of microbial genomes recovered from isolates, single cells, and
853        metagenomes. *Genome Res* 2015. https://doi.org/10.1101/gr.186072.114.

22

854  74  Ondov BD, Treangen TJ, Melsted P, Mallonee AB, Bergman NH, Koren S, *et al.* Mash:
855      fast genome and metagenome distance estimation using MinHash. *Genome Biol*
856      2016;**17**:132. https://doi.org/10.1186/s13059-016-0997-x.
857  75  Li H. A statistical framework for SNP calling, mutation discovery, association mapping
858      and population genetical parameter estimation from sequencing data. *Bioinformatics*
859      2011;**27**:2987–93. https://doi.org/10.1093/bioinformatics/btr509.
860  76  Jain C, Rodriguez-R LM, Phillippy AM, Konstantinidis KT, Aluru S. High throughput ANI
861      analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat Commun*
862      2018;**9**:5114. https://doi.org/10.1038/s41467-018-07641-9.
863  77  Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using DIAMOND.
864      *Nat Methods* 2015;**12**:59–60. https://doi.org/10.1038/nmeth.3176.
865  78  Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, *et al.* SNP-sites:
866      rapid efficient extraction of SNPs from multi-FASTA alignments. *Microb Genom*
867      2016;**2**:e000056. https://doi.org/10.1099/mgen.0.000056.
868  79  Price MN, Dehal PS, Arkin AP. FastTree 2--approximately maximum-likelihood trees for
869      large alignments. *PLoS One* 2010;**5**:e9490.
870  80  Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y. Ggtree: An r package for visualization and
871      annotation of phylogenetic trees with their covariates and other associated data.
872      *Methods Ecol Evol* 2017;**8**:28–36. https://doi.org/10.1111/2041-210x.12628.
873  81  Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
874      stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol*
875      2015;**32**:268–74. https://doi.org/10.1093/molbev/msu300.
876  82  Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool.
877      *J Mol Biol* 1990;**215**:403–10. https://doi.org/10.1016/S0022-2836(05)80360-2.
878  83  Kassambara A. *rstatix: Pipe-friendly Framework for Basic Statistical Tests in R*. 2023.
879      URL: https://github.com/kassambara/rstatix (Accessed 19 December 2023).
880  84  Wickham H. *Ggplot2*. New York, NY: Springer; 2011.
881  85  *drawio: draw.io is a JavaScript, client-side editor for general diagramming and*
882      *whiteboarding*. 2023. URL: https://github.com/jgraph/drawio (Accessed 25 September
883      2023).
884  86  Bogart S. *sankeymatic: Make Beautiful Flow Diagrams*. 2023. URL:
885      https://github.com/nowthis/sankeymatic (Accessed 25 September 2023).
886  87  Jalil M, Quddos F, Anwer F, Nasir S, Rahman A, Alharbi M, *et al.* Comparative Pan-
887      Genomic Analysis Revealed an Improved Multi-Locus Sequence Typing Scheme for
888      Staphylococcus aureus. *Genes* 2022;**13.**: https://doi.org/10.3390/genes13112160.
889  88  Liu N, Liu D, Li K, Hu S, He Z. Pan-Genome Analysis of Staphylococcus aureus
890      Reveals Key Factors Influencing Genomic Plasticity. *Microbiol Spectr*
891      2022;**10**:e0311722. https://doi.org/10.1128/spectrum.03117-22.
892  89  Bosi E, Monk JM, Aziz RK, Fondi M, Nizet V, Palsson BØ. Comparative genome-scale
893      modelling of Staphylococcus aureus strains identifies strain-specific metabolic
894      capabilities linked to pathogenicity. *Proc Natl Acad Sci U S A* 2016.
895      https://doi.org/10.1073/pnas.1523199113.
896  90  Park S, Jung D, O'Brien B, Ruffini J, Dussault F, Dube-Duquette A, *et al.* Comparative
897      genomic analysis of Staphylococcus aureus isolates associated with either bovine
898      intramammary infections or human infections demonstrates the importance of
899      restriction-modification systems in host adaptation. *Microb Genom* 2022;**8.**:
900      https://doi.org/10.1099/mgen.0.000779.
901  91  Sassi M, Bronsard J, Pascreau G, Emily M, Donnio P-Y, Revest M, *et al.* Forecasting
902      Staphylococcus aureus Infections Using Genome-Wide Association Studies, Machine
903      Learning, and Transcriptomic Approaches. *mSystems* 2022:e0037822.
904      https://doi.org/10.1128/msystems.00378-22.
905  92  Aanensen DM, Feil EJ, Holden MTG, Dordel J, Yeats CA, Fedosejev A, *et al.* Whole-
906      Genome Sequencing for Routine Pathogen Surveillance in Public Health: a Population
907      Snapshot of Invasive Staphylococcus aureus in Europe. *MBio* 2016;**7.**:
908      https://doi.org/10.1128/mBio.00444-16.
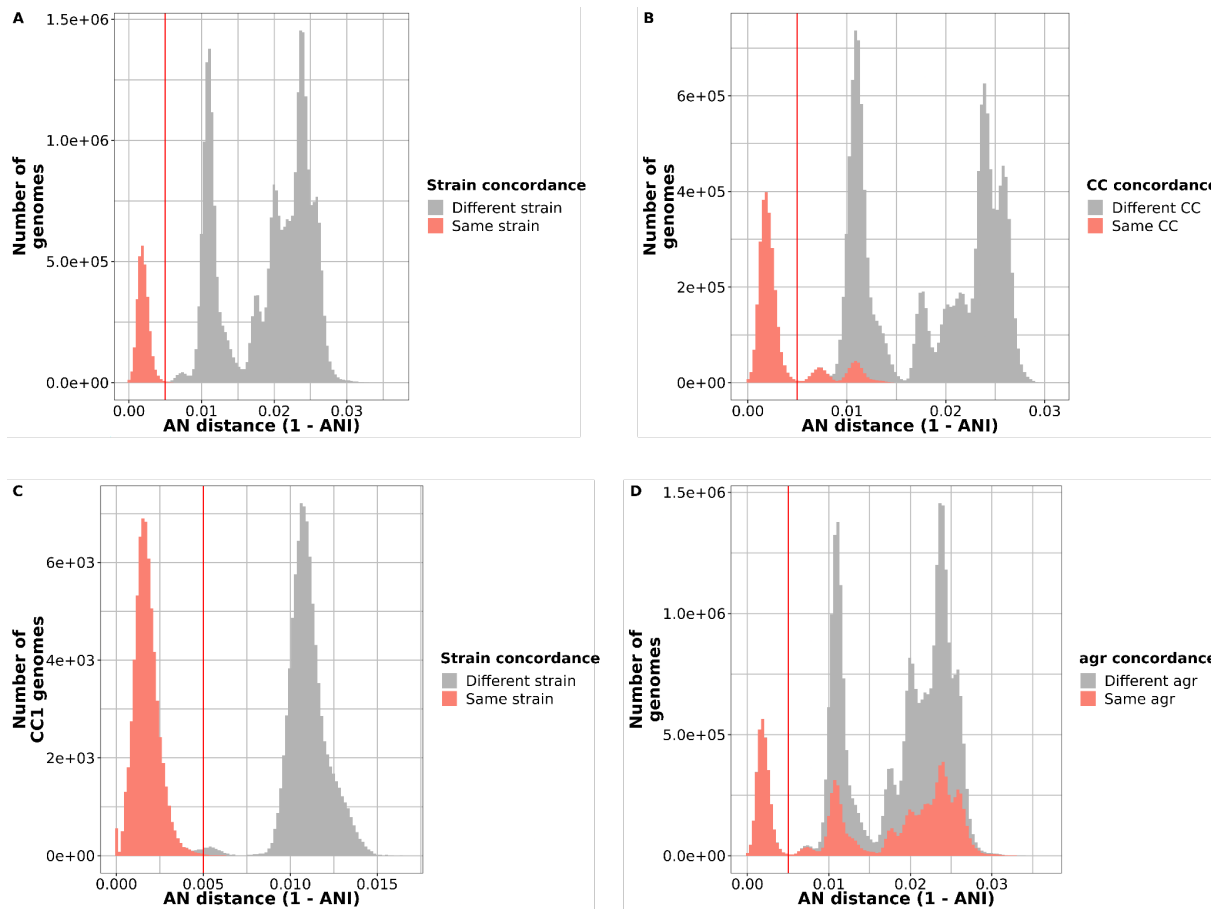
23

909 <u>Figures</u>

910 **Figure 1**



911

912 Figure 1: Sankey diagram showing the fate of 83,383 *S. aureus* whole genome shotgun
913 datasets and 1475 complete genomes through processing and filtering.
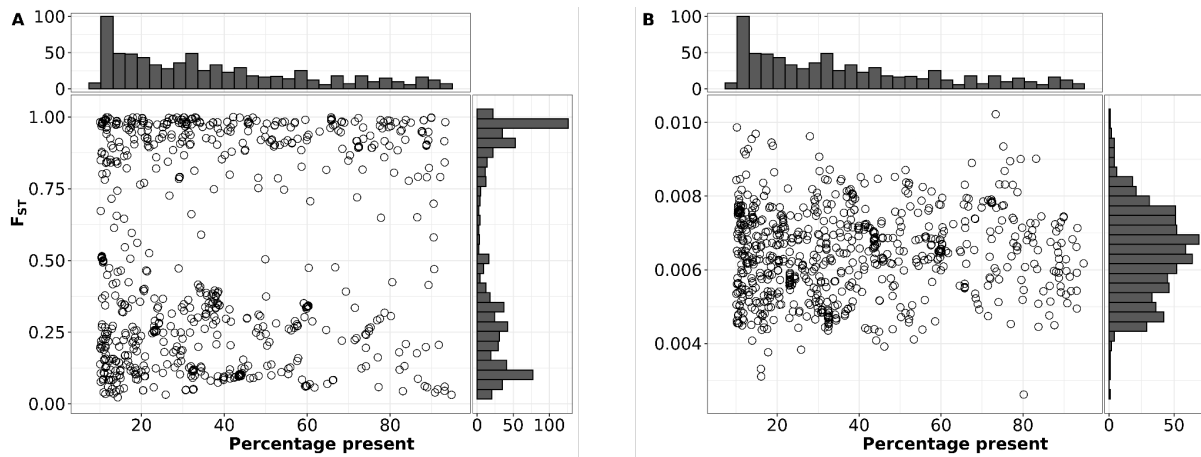
914

24

**Figure 2**



Figure 2: An average nucleotide identity of >99.5% in the core genome defines the strain boundary of *S. aureus*.
For our dataset of 7954 substrains, all-vs-all pairwise Average Nucleotide (AN) distances were plotted as a histogram. (A) Sample pairs less than 0.005 AN distance apart (i.e. greater than 99.5% ANI) were grouped as a strain. ( B) Strains and clonal complex designations don't exactly overlap. The pairwise AN distance histogram was colored by whether the genomes were in the same clonal complex. (C) CC1 genomes are in different strains. AN distances of genomes assigned to CC1 showing that there are within- and between- strain distances. (D) Genomes in the same strain have the same *agr* group.The pairwise ANI distance histogram was colored by whether the genomes were in the same *agr* group.

25

**Figure 3**



Figure 3: Bimodal distribution of $F_{ST}$ for intermediate genes.

Each circle represents an individual intermediate gene from the 7954 substrain pangenome. Percentage prevalence on the x axis is the percentage of genomes the gene is found in. $F_{ST}$ or 'fixation index' is on the y axis. (A) $F_{ST}$ scores calculated for each intermediate gene with 99.5% ANI-based clustering. (B) As a control, $F_{ST}$ scores were calculated for each intermediate gene when clusters were randomly assigned.

26

937 **Figure 4**



938

939 Figure 4: Strain-group specificity and co-occurrence of specific Staphylococcal toxins.
940 Core genome of the 7954-set. Heatmap on right shows presence absence and $F_{ST}$ of specific
941 Staphylococcal toxins - Panton-Valentine Leukocidin (LukF and LukS), Toxic Shock Syndrome
942 Toxin (TSST), and Staphylococcal Enterotoxins type A, B, G, U (SEA, SEB, SEG, SEU),
943 Superantigen like protein (SSL8), Leukocidin ED (LukE, LukD) . The colors of the whole-genome
944 phylogeny are based on strain assignments.

945

27

946 **Figure 5**



947

948 Figure 5: Relationship between gene prevalence, number of strains and homoplasy for
949 non-core genes.
950 Each dot represents a non-core gene in the 740-set pangenome. Purple = rare genes, green =
951 concentrated, Brown = diffuse. In panel a, the relationship between overall prevalence (number of
952 genomes out of 740) and number of strains (out of 37) each gene is found in is shown. The curves
953 for the theoretical minimum number of strains for a given number of genomes (x/20) is shown in
954 solid black and the extreme random distribution (37*(1-exp(-x/37))) is shown in dashed black. Panel
955 b shows the relationship between prevalence of estimated number of changes on the species tree
956 calculated by homoplasyfinder[27].

957

28

958 **Figure 6**



959

960 Figure 6: t-SNE analysis of 740-seq differentiated by non-core gene sets
961 Each dot represents one of the genomes of the 740-set colored by its strain membership. Different
962 sets of non-core genes were used as input for the t-SNE: a) only rare; b) only strain-diffuse; c) only
963 strain-concentrated; d) all non-core.
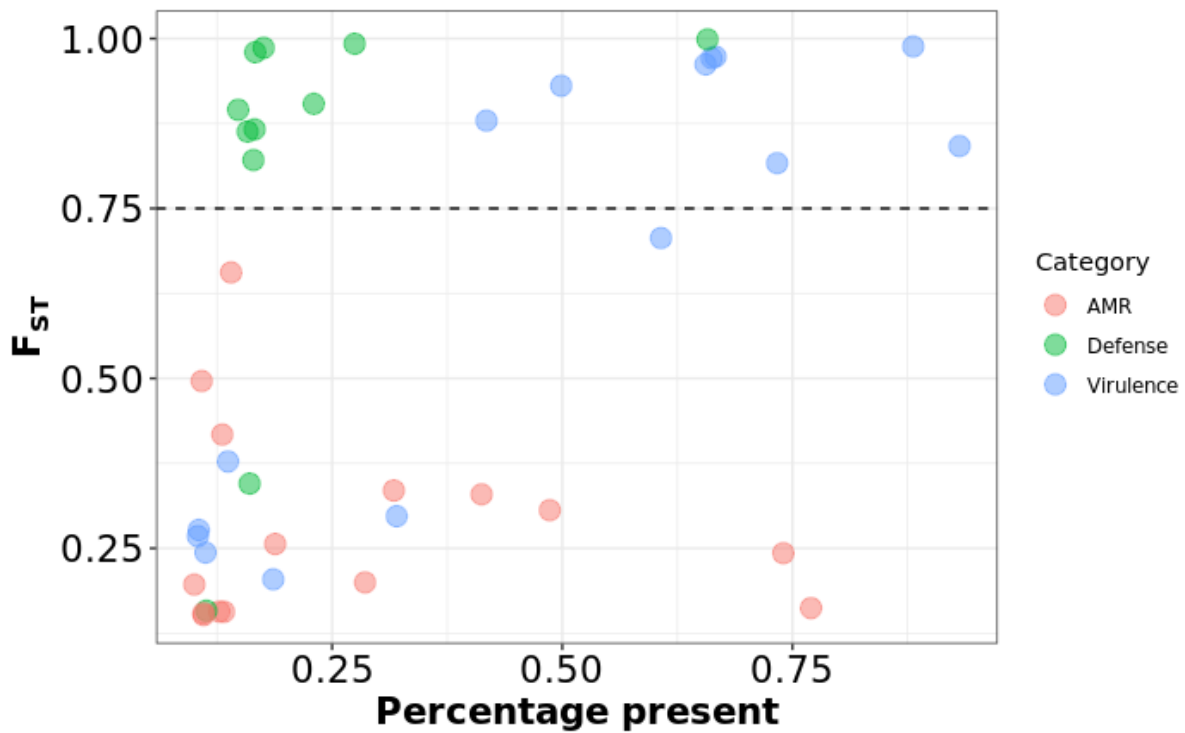
964

29

965    **Figure 7**



966



967

968    Figure 7: Distribution of different categories on non-core genes on the *S. aureus*
969    chromosome using two orthologous methods.
970    A: Location based on 337 complete genome sequences. The start site for every gene in each
971    category was obtained for 337 chromosomes. The totals were placed in 10,000 bp bins on the
972    chromosome and the proportion of the total for each class is plotted (i.e. the sum of the values of
973    the 10,000 bins is 1). Purple = rare genes; green = strain-concentrated; brown = strain-diffuse. B:
974    Location based on the nearest core gene. For all 7,954 substrains, the closest core gene on the
975    same contig was determined. The x axis are start sites for the core genes of genome N315
976    (GCA_000009645)[29]. The values were binned and proportionalized as in A. For both A and B the
977    location of selected features is shown: I = SCCmec; II = type VII secretion system; III = vSaα; IV =
978    phiSa1; V=vSaγ; VI = phiSa2; VII = vSaβ; VIII = phiSa3; IX=vSa4). N315 coordinates are based on
979    Gill et al [29] and Warne et al [30], except phiSa2 and phiSa3, which are from Mu50 and MW2,
980    respectively.

30

981    **Figure 8**
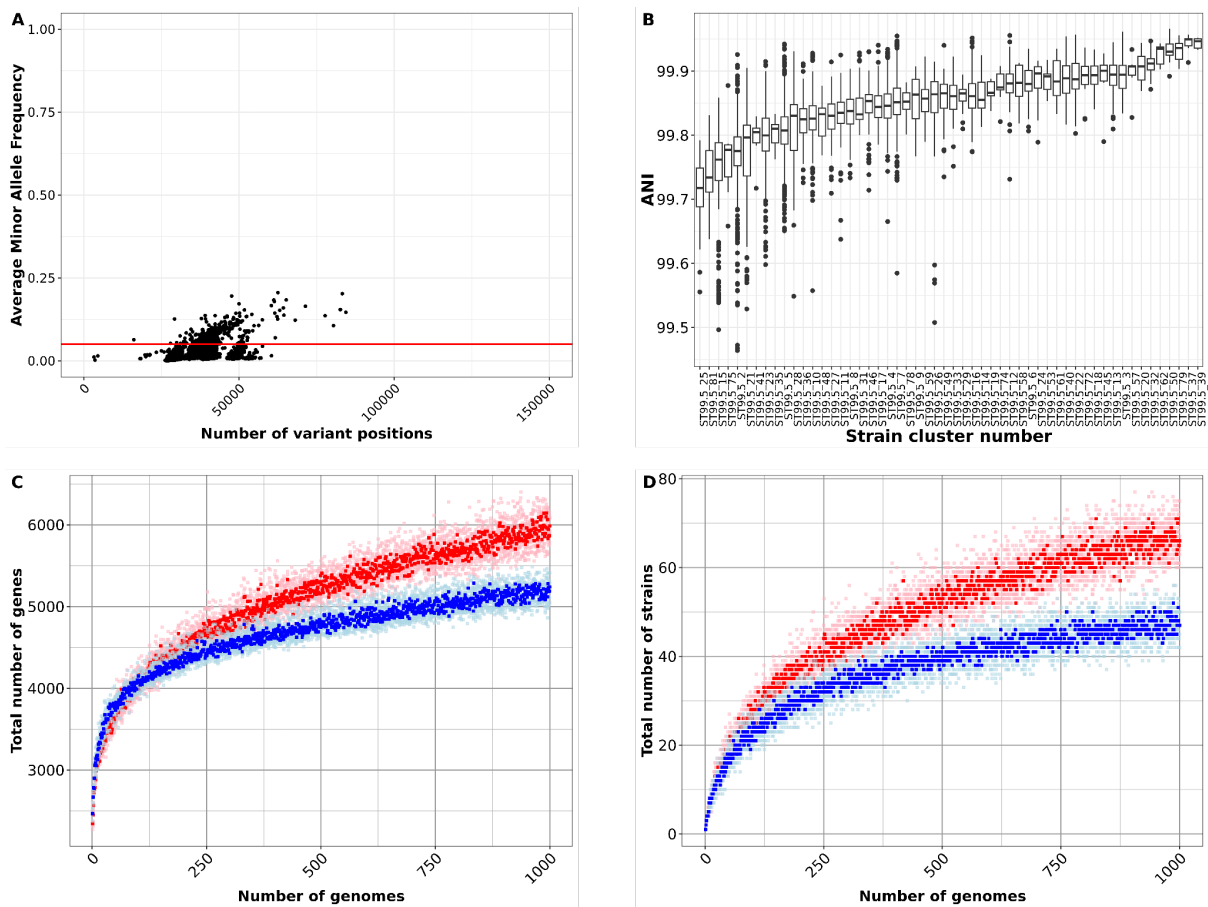


982

983    Figure 8: Prevalence versus Fst for intermediate antimicrobial-resistance (AMR),
984    virulence and phage defense genes
985    AMR and virulence genes were identified using AMRFinder+[31], phage defense genes were
986    identified using defense-finder[33]. The dashed horizontal line represents the boundary between
987    strain-diffuse and strain-concentrated.
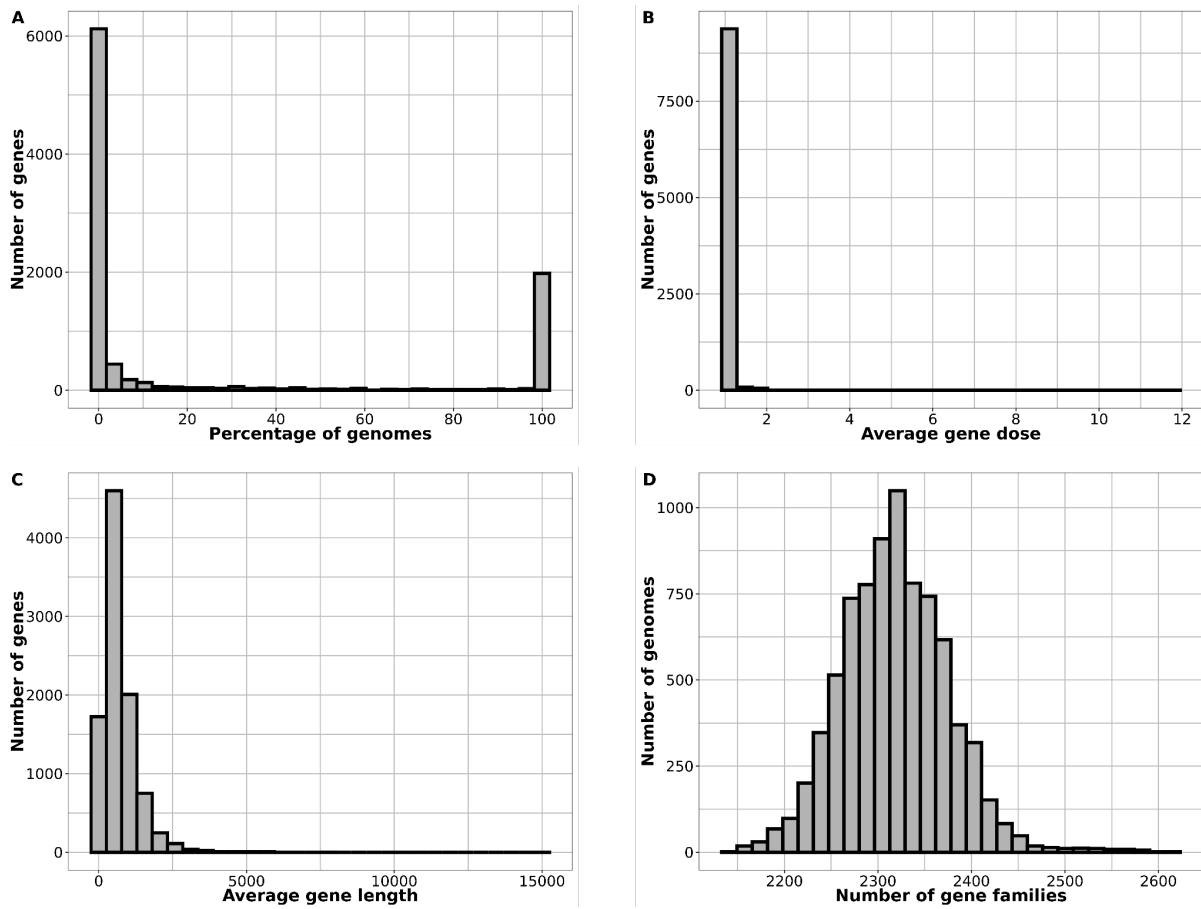
31

988    Supplemental Data



989

990    <u>Figure S1: Effect of filtering, clustering and dereplicating 83,383 *S. aureus* genomes</u>
991    (A) The x-axis shows the total number of variants when compared with the Bactopia auto-chosen
992    reference, and the y-axis shows the average minor allele frequency (MAF). Each dot is one of
993    57,093 genomes which were obtained after filtering out samples ranked 'Bronze' or 'Exclude' by
994    Bactopia and/or found to have non-*S. aureus* genome content by Bactopia and CheckM (Figure 1).
995    Samples in the top quadrant (Above red horizontal line - Average MAF > 0.05) were considered to
996    be *S. aureus* strain mixtures and were discarded. The remaining 56,771 samples in the bottom
997    quadrant (< 0.05 Average MAF) were used for further analysis. (B) Boxplots showing spread of
998    pairwise ANI within each "strain" cluster. Only strain clusters having more than 10 genomes are
999    shown. Black horizontal line within each boxplot shows the median within strain-cluster pairwise
1000   ANI. Total number of unique genes discovered (C) and total number of strains discovered (D) for
1001   every new genome added from the dereplicated set (red dots) or a random genome from the un-
1002   dereplicated 58,034 (blue dots). Up to 1000 random genomes were added from each set and the
1003   total number of unique genes or strains were measured for every genome added (light red and
1004   light blue dots). This procedure was repeated 5 times and the median number of genes or strains
1005   discovered are shown in dark red and dark blue dots. More genes and more strains were
1006   discovered from the same number of genomes (after observing 1000 genomes) in the dereplicated
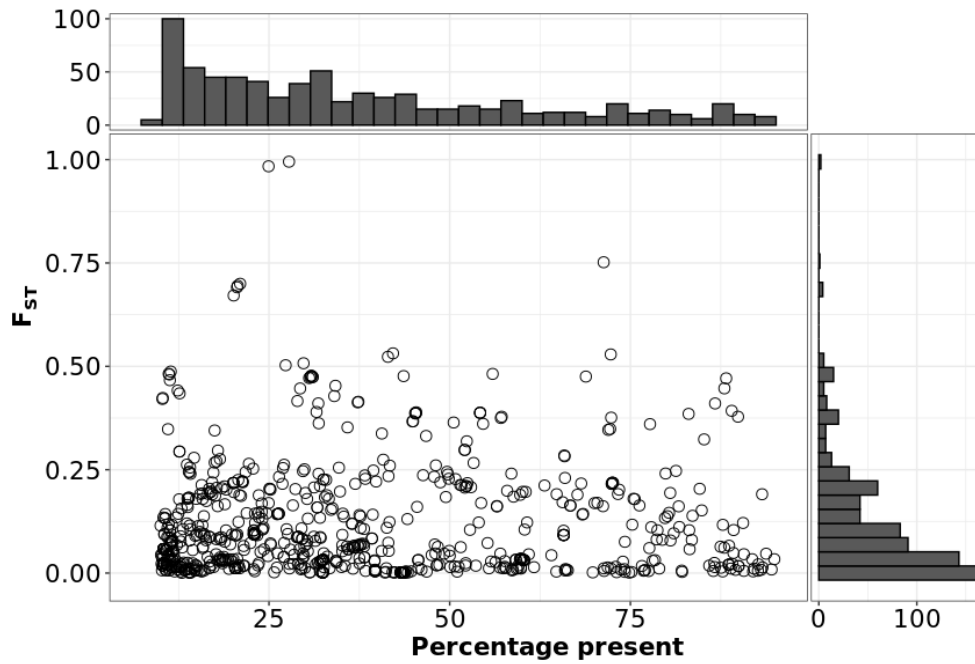1007   set compared to the un-dereplicated set.

32

1008



1009

Figure S2: The 7,954 substrain pangenome of *S. aureus*.
Histograms depicting the (A) frequency distribution of genes in our dataset, (B) the average dosage of each gene per genome, (C) the average length distribution of each gene, and (D) the distribution of the number of unique PIRATE gene families per genome.
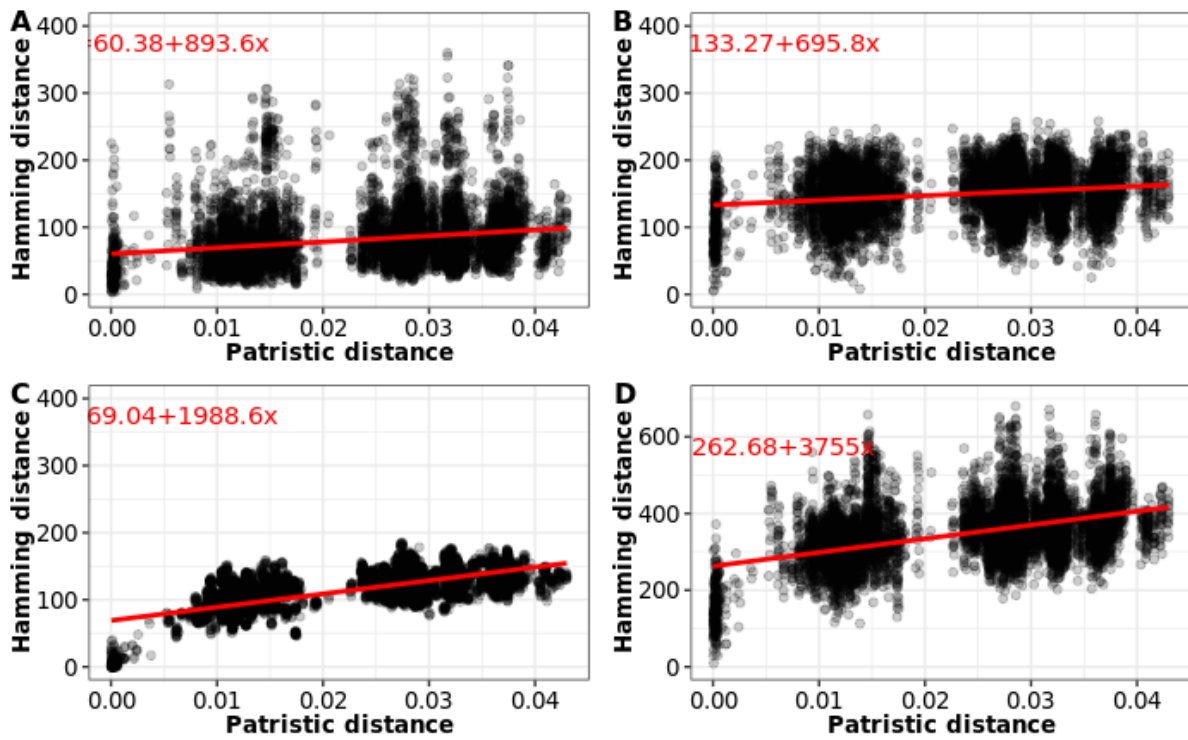
33

1014



1015

1016  Figure S3: There are no *agr* group specific intermediate genes aside from *agrD*.
1017  Dot plot showing percentage prevalence of only intermediate genes (> 10%, < 95%) on the x-axis
1018  and the corresponding $F_{ST}$ on the y-axis. $F_{ST}$ scores calculated for *agr* type-based population
1019  segregation. The three dots > 0.75 $F_{ST}$ correspond to the *agrD*, which are known to be lineage
1020  specific . The *agrD* of the fourth *agr* type is absent in this plot as it is present in < 10% of the
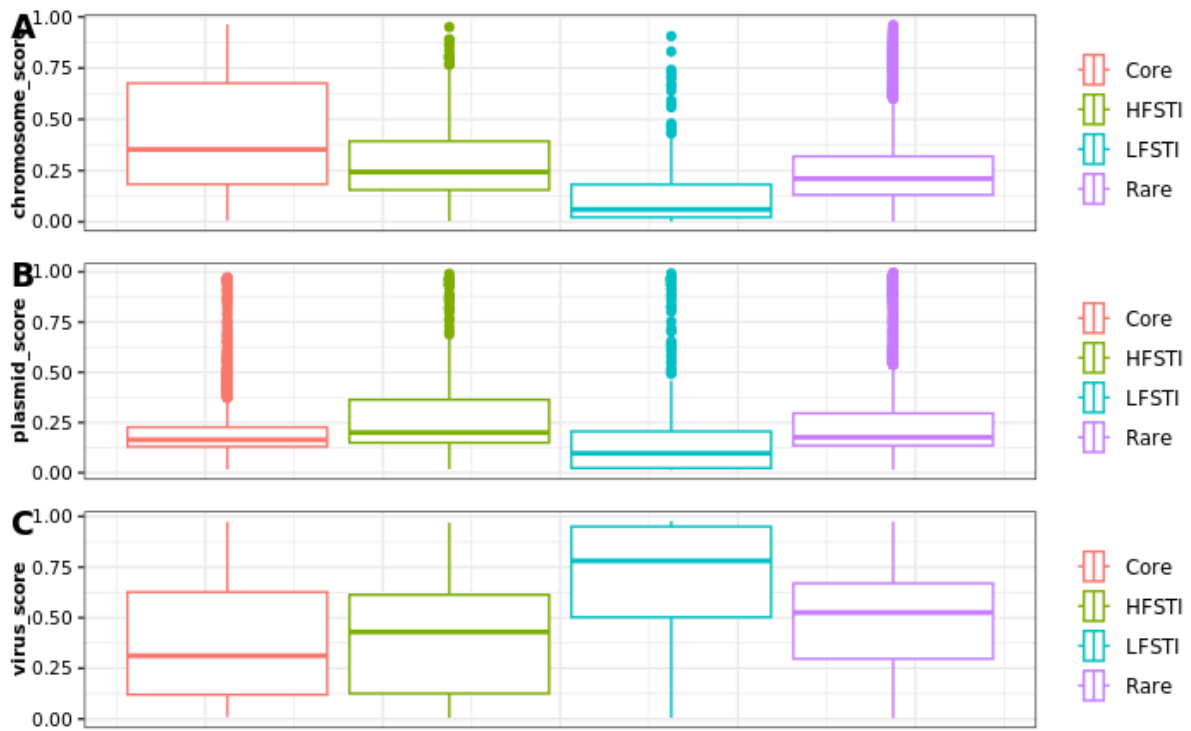1021  population [15].

34

1022



1023

1024 Figure S4: strain-concentrated gene content declines gradually with core-gene distance.
1025 Each dot represents a comparison between substrains in the 740-set. Patristic distance was tip-tip
1026 distance on the phylogeny. Hamming distance was calculated from a presence absence matrix of
1027 each non-core gene type: A) rare genes, B) strain-diffuse, C) strain-concentrated, D) all non-core
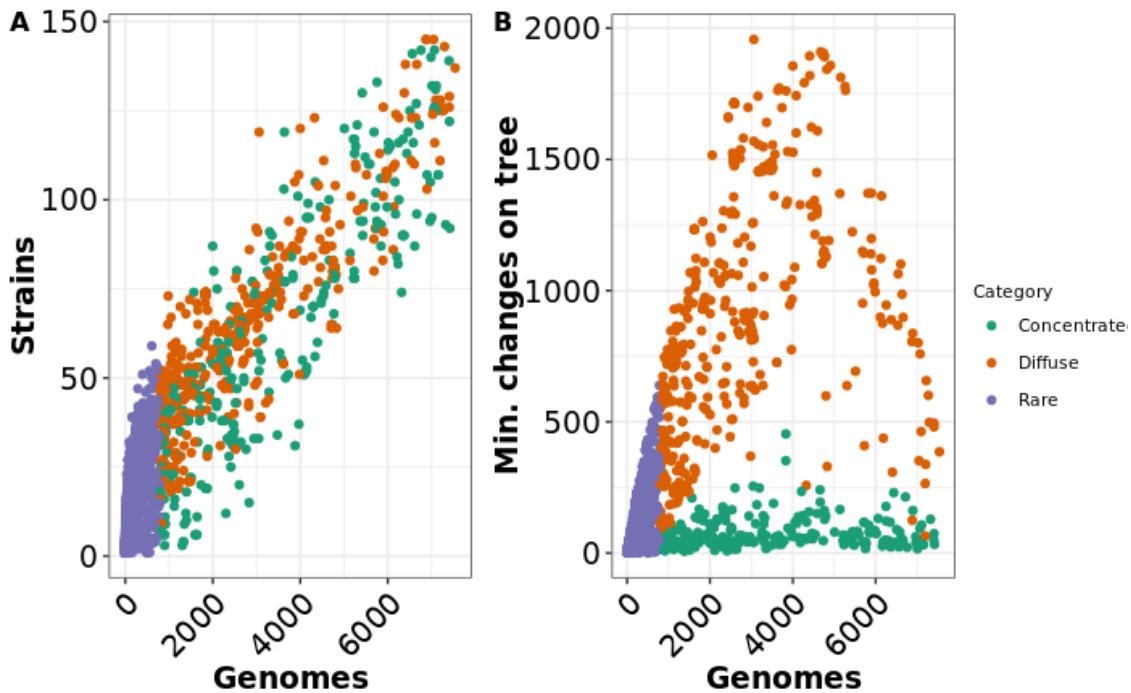1028 (note different y-axis scale). Red lines show the linear model fit.

35

1029



1030

Figure S5: Genomad score distributions for 7954-set pangenomes.
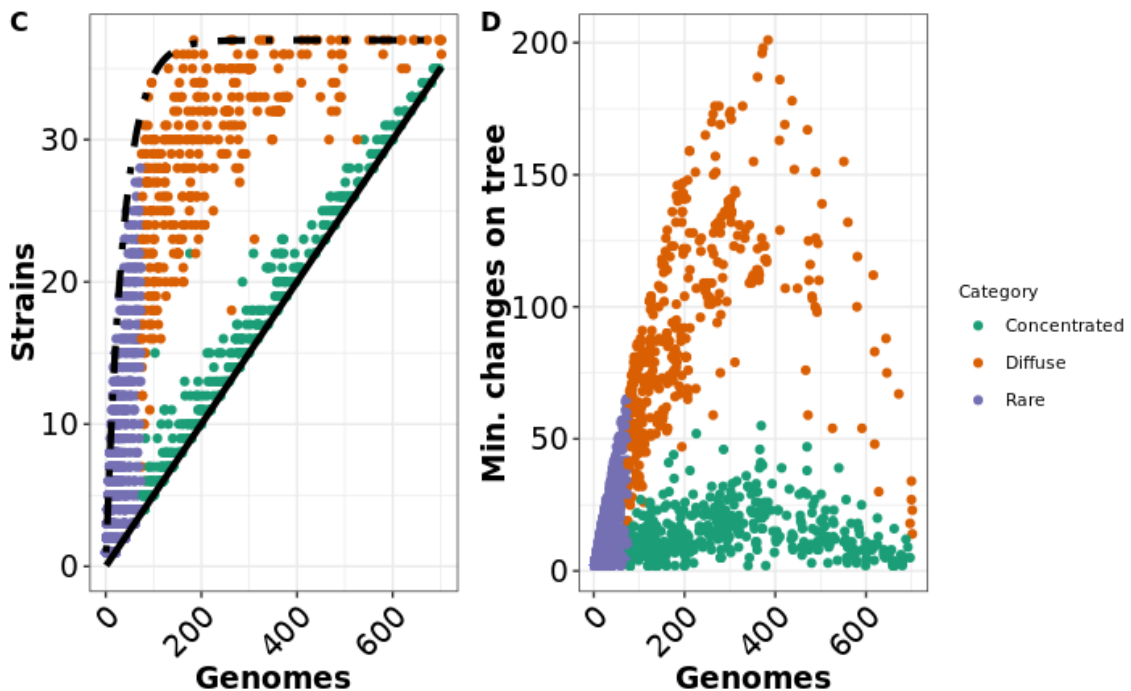The geNomad [28] probability scores for A) chromosome B) plasmid and C) virus were grouped by
gene class. All differences were significant in a Tukey's pairwise comparisons at < 0.05 (corrected
for multiple tests), except strain-diffuse-Core plasmid_score and strain-concentrated-Core virus
score.

36

1036
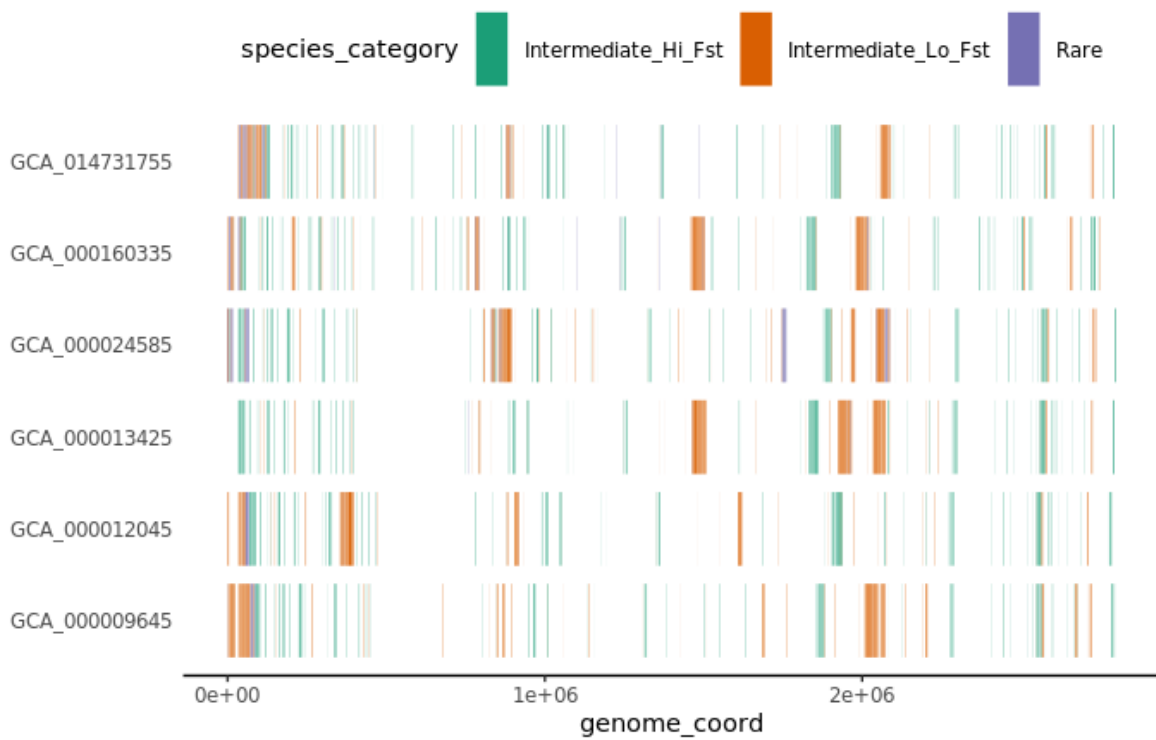


1037

1038

Figure S6 Relationship between gene prevalence, number of strains and homoplasy for non-core genes for the 7954-set (a,b) and the 740-set-90 (c,d)

The plots are formatted as Figure 5. Each dot represents a non-core gene. "concentrated" = strain-concentrated, "diffuse" = strain-diffuse. Panels A and C show the relationship between overall prevalence (number of genomes out of 740) and number of strains (out of 37) each gene is found in. Panels B and D show the relationship between prevalence of estimated number of changes on the species tree calculated by homoplasyfinder[27]. In panel A, the unbalanced nature of the 7954-set (a few strains have thousands of genomes, many have only one) obscures the differences between concentrated and diffuse: it not possible to plot simple bounds of lowest possible and random gene distribution into strains as it is for the 740-90 set (panel C).

37

1049



Figure S7: Chromosome start locations of non-core genes on six S.aureus complete chromosomes.
The name on the left-hand side refers to NCBI assembly database designations. GCA_014731755 is CC30 MRSA; GCA_000160335 is CC30 MSSA; GCA_000024585 is CC5 MSSA; GCA_0000134525 is CC8 MRSA; GCA_000012045 is CC8 MRSA; GCA_000009645 is CC5 MRSA (N315 the S. aureus type strain). "Intermediate_Hi_Fst" = strain-concentrated; "Intermediate_Lo_Fst" = strain-diffuse.

38

1058   Table S1: *S. aureus* studies quoting pangenome statistics.
1059   "?" indicates that the corresponding information could not be found

| Title | Date | No. of genomes | Sampling space | Assembly level | Pangenome tool | No. core | Total gene families |
|---|---|---|---|---|---|---|---|
| Comparative Pan-Genomic Analysis Revealed an Improved Multi-Locus Sequence Typing Scheme for Staphylococcus aureus [87] | 2022-11-19 | 502 | Diverse | Complete | PanRV (Roary) | 2320 | 12477 |
| Pan-Genome Analysis of Staphylococcus aureus Reveals Key Factors Influencing Genomic Plasticity [88] | 2022-11-01 | 1519 | Diverse | All | Roary | 1000 | 16794 |
| Pangenomic Approach To Understanding Microbial Adaptations within a Model Built Environment, the International Space Station, Relative to Human Hosts and Soil [50] | 2022-01-08 | 106 | ISS, human, soil | All | Roary | 1935 | 6847 |
| The Epidemiological and Pangenome Landscape of Staphylococcus aureus and Identification of Conserved Novel Candidate Vaccine Antigens [47] | 2022-02-01 | 355 | Diverse | All | ? | 2025 | 7199 |
| Analyses of Livestock-Associated Staphylococcus aureus Pan-Genomes Suggest Virulence Is Not Primary Interest in Evolution of Its Genome [51] | 2019-05-22 | 14 | Livestock associated | Complete | Roary | 1969 | 4637 |
| Comparative genome-scale modelling of Staphylococcus aureus strains identifies strain-specific metabolic capabilities linked to pathogenicity [89] | 2016-06-10 | 64 | Diverse | All | dGenome DuctAPE | 1441 | 7457 |
| PanRV: Pangenome-reverse vaccinology approach for identifications of potential vaccine candidates in microbial pangenome [48] | 2019-03-12 | 301 | Diverse | All | PanRV (Roary) | 1524 | 11384 |
| Whole-Genome Sequencing of Staphylococcus aureus and Staphylococcus haemolyticus Clinical Isolates from Egypt [43] | 2022-06-21 | 90 | 56 from 1 hospital and 34 from greater Arab region | All | Anvio | 1501 | 4283 |
| Phylogenomic Analysis Reveals the Evolutionary Route of Resistant Genes in Staphylococcus aureus [52] | 2019-11-03 | 152 | Diverse | Complete | Manual alignment and clustering | 2426 | 6326 |
| Comparative genomic analysis of Staphylococcus aureus isolates associated with either bovine intramammary infections or human infections demonstrates the importance of restriction-modification systems in host adaptation [90] | 2022-02-18 | 187 | Human and cattle | All | Roary | 2700 | 6812 |

39

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Molecular Epidemiology of Staphylococcus aureus in China Reveals the Key Gene Features Involved in Epidemic Transmission and Adaptive Evolution [44] | 2022-10-03 | 332 | Human clinical strains from China | All | Heap's law algorithms | 890 | 5832 |
| Estimated Roles of the Carrier and the Bacterial Strain When Methicillin-Resistant Staphylococcus aureus Decolonization Fails: a Case-Control Study [45] | 2022-08-24 | 477 | MRSA carriers from Denmark hospitals | All | panX | 1671 | 5925 |
| Forecasting Staphylococcus aureus Infections Using Genome-Wide Association Studies, Machine Learning, and Transcriptomic Approaches [91] | 2022-07-05 | 356 | Mostly human | All | Panaroo | 1489 | 8827 |
| Carriage prevalence and genomic epidemiology of Staphylococcus aureus among Native American children and adults in the Southwestern USA [46] | 2022-05-13 | 92 | Native Americans from Southwestern USA | Complete | Roary | 1808 | ? |
| Polyclonality, Shared Strains, and Convergent Evolution in Chronic Cystic Fibrosis Staphylococcus aureus Airway Infection [42] | 2020-03-23 | 1382 | Longitudinal sampling from 246 children with CF from the US | All | Roary | 1142 | 21358 |
| PIRATE: A fast and scalable pangenomics toolbox for clustering diverged orthologues in bacteria [17] | 2019-10-09 | 253 | Diverse | All | PIRATE | 2433 | 4250 |
| Whole-Genome Sequencing for Routine Pathogen Surveillance in Public Health: a Population Snapshot of Invasive Staphylococcus aureus in Europe [92] | 2016-05-05 | 308 | Invasive isolates from Europe hospitals within a 6 month period | All | BlastP & TribeMCL | ? | 4281 |

1060