1

# Pan-genome analysis reveals the molecular basis of niche adaptation of *Staphylococcus epidermidis* strains

Fei Su[a,b*], Rui Tian[c*], Yi Yang[d*], Lihui Zou[b], Xiaomao Xu[e], Dongke Chen[f], Junhua Zhang[b], Xue Chen[f], Fei Xiao[a,b], Gang Zhao[g], Yanming Li[e#], Hongtao Xu[f#]

[a]Clinical Biobank, [b]The Key Laboratory of Geriatrics, [d]Department of Otorhinolaryngology [e]Department of Respiratory and Critical Care Medicine, [f]Department of Laboratory Medicine, [g]Department of Surgery, Beijing Hospital, National Center of Gerontology, Beijing 100730, P. R. China

[c]Department of Cardiovascular Disease, Beijing Anzhen Hospital, Capital Medical University, Beijing 100024, People's Republic of China

*These authors contributed equally to this work

#Address correspondence to: Hongtao Xu, xuhongtao2911@bjhmoh.cn or Yanming Li, liyanming2632@bjhmoh.cn

**Running title:** Pan-genome analysis of Staphylococci

**Keywords:** *Staphylococcus epidermidis*, antimicrobial resistance, pan-genome, mobile genetic element

## Abstract

*Staphylococcus epidermidis* is the most commonly isolated species from human skin and the second leading cause of bloodstream infections. Here, we performed a large-scale comparative study without any pre-assigned reference to identify genomic determinants associated with their diversity and adaptation as a "double-side spy", a skin dominant colonization, and a successful pathogen. The pan-genome of *S. epidermidis* is open with 435 core proteins and a pan-genome size of 8034 proteins. Genome-wide phylogenetic tree shows that whole genome sequence is a powerful tool to analyze the complex evolutionary process of *S. epidermidis* and investigate the source of infection. Comparative genome analyses demonstrate the high diversity of antimicrobial resistances, especially mobile genetic elements. The complicated relationships of host-bacterium and bacterium-bacterium help *S. epidermidis* to play a vital role in balancing the epithelial microflora. The highly variable and dynamic nature of the *S. epidermidis* genome may be the result of its success in adapting to broad habitats, which is necessary to deal with complex environments. This study gives the general landscape of *S. epidermidis* pan-genome and provides valuable insights into mechanisms for genome evolution and lifestyle adaptation of this ecologically flexible species.

## Introduction

The coagulase-negative *Staphylococcus epidermidis* (*S. epidermidis*) is a common human skin commensal bacterium that can be cultured from virtually every body surface of healthy individuals. It also plays a central role in the skin microbiome [1, 2], it can keep the ecological balance of human skin microflora [3]. *S. epidermidis* can produce various bacteriocins, which kill other microorganisms and have frequently been proposed to enhance survival of the producer strains in a competitive fashion [4, 5]. Especially, serine protease Esp, secreted by *S. epidermidis*, can inhibit the biofilm formation of *S. aureus* and destroy pre-existing *S. aureus* biofilms [6].

However, *S. epidermidis* is the second most common cause of nosocomial infections, which in most cases are antibiotic-resistant [1, 7]. Antibiotic resistance remarkably complicates the treatment and increases the medical expenses [8]. The large gene pool of antibiotic resistance in *S. epidermidis* is shared with many other pathogenic species such as *S. aureus* [9]. Mobile genetic element, multidrug-resistant conjugative plasmids, arginine catabolic mobile element (ACME) [9], and staphylococcal chromosome cassette *mec* (SCC*mec*) elements [10] conferring β-lactam resistance are transferred frequently, enabling rapid evolution and adaptation against antibiotic selection pressure [11, 12]. When the protective layer of the human epithelium is breached and the mechanisms of host immunity fail, staphylococcal infections can become extremely dangerous and even fatal [13]. *S. epidermidis* is particularly associated with the increased use of indwelling medical devices such as artificial heart valves, prosthetic joints, and vascular catheters, which provide a substrate for biofilm formation. On the other hand, during the long-time "arms race", human beings have developed versatile immunity system with antimicrobial peptides (AMPs) as the first line of innate immune defense on the human skin; meanwhile, *S. epidermidis* also owns multiple mechanisms such as surface charge alteration, extracellular proteases, exopolymers, and efflux pump proteins

70  to fight against AMPs [7]. The complex host-bacterium and bacterium-bacterium

71  relationships make it necessary to investigate the genetic diversity, genome evolution, and

72  lifestyle adaptation of *S. epidermidis*.

73  Much attention has been focused on understanding the evolution and spread of *S. epidermidis*

74  by different methods [11, 14]. As the time goes on, high throughput sequencing is now fast

75  and cheap and a large amount of genomics data about *S. epidermidis* are accumulated, it is

76  essential to perform more comprehensive comparative and evolutionary study of ecologically

77  diverse strains of *S. epidermidis* for better clinical management. Here, we compared the

78  genomic features of *S. epidermidis* isolates of clinical and non-clinical relevance by using a

79  pan-genome analysis of 198 publicly available *S. epidermidis* strains at the GenBank

80  database of National Center for Biotechnology Information at April 30 of 2017. We

81  assembled the consensus "pan-chromosome" without any pre-assigned genome reference and

82  identified both core and variable regions within the chromosome. Second, we utilized a

83  comparative genomics approach on 198 genomes to analyze the diversity of antibiotic

84  resistance of *S. epidermidis*. Our results revealed that *S. epidermidis* isolates encoded a vast

85  collection of genetic determinants and mechanisms to confer antibiotic resistance,

86  antimicrobial peptides resistance, and survival adaptations. These analyses will provide

87  insight into the coevolution of *S. epidermidis* as a nosocomial pathogen and directly aid the

88  future efforts for large-scale epidemiological studies of this continuously evolving multi-drug

89  resistant organism.

90

## Methods

### Strains

A total of 198 *S. epidermidis* isolates were selected to represent known diversity within the species and multiple locations and sources until Apirl 30 of 2017, including reference genomes from S. epidermidis strain RP62A (Gill et al. 2005). All the available genome sequence of *S. epidermidis* strains and related annotation data were downloaded through the GenBank database [15] of NCBI (see Table S1 in the supplemental material).

### SCC*mec* and ACME typing

An SCC*mec* sequence cassette database was prepared with the following accession numbers downloaded from NCBI: AB033763.2 (Type I), AB433542.1 (Type I.2), D86934.2 (Type II), AB261975.1 (Type II.4), AJ810123 (Type II-B), AB127982.1 (Type II-B), AM983545.1 (Type II-D), HE858191.1 (Type II-E), AB037671.1 (Type III), HM030721.1 (Type IV), HM030720.1 (Type IV), AM292304.1 (ZH47 mobile elements), AB425824.1 (Type IV), EU437549.2 (Type IV-A), AB063172.2 (Type IV-A), AB063173 (Type IV-B), AY271717.1 (Type IV-C), AB096217 (Type IV-C), AB245470.1 (Type IV-C), AB097677.1 (Type IV-D), AJ810121.1 (Type IV-E), DQ106887.1 (Type IV-G), AB633329.1 (Type IV-I), AB425823.1 Type IV), AB121219.1 (Type V), AB478780.1 (Type V), AB512767.1 (Type V), AF411935.3 (Type VI), AB462393.1 (Type VII), AB373032.1 (Type V-C1), FJ670542.1 (Type VIII), FJ390057.1 (Type VIII), AB505628.1 (Type IX), AB505630.1 (Type X), and FR821779.1 (Type XI) [16].

The ACME-*arc*A and ACME-*opp*3AB genes were used as markers of the ACME-*arc* cluster and the ACME-*opp*3 cluster, respectively. ACME was classified as type I (contains the ACME-*arc*A and ACME-*opp*3AB gene clusters), type II (carries only the ACME-*arc*A locus), and type III (carries only the ACME-*opp*3AB locus) [17]. ACME-*arc*A and ACME-

115  *opp*3AB identified in this study were compared with the reference sequences of ACME-*arc*A

116  (USA300_FPR3757) and ACME-*opp*3AB (USA300_FPR3757).

117  **kSNP *S. epidermidis* trees.**

118  A phylogenetic tree was inferred from single-nucleotide polymorphisms (SNPs) identified by

119  kSNP (version: 3.0, https://sourceforge.net/projects/ksnp/) [18] by using a *k*-mer length 19

120  nucleotides and based on a requirement that at least 80% of the genomes have a nucleotide at

121  a given SNP position in order for the SNP to be considered to be a core and included in tree

122  building. A total of 1832 core SNP positions were identified. These SNPs were used to infer

123  a maximum-likelihood tree with RAxML [19] with 100 bootstrap replicates.

124  **Pan-genome analysis**

125  Cluster of orthologous proteins were generated with version 3.24 of PanOCT

126  (https://sourceforge.net/projects/panoct/) as previous described [20]. Briefly, PanOCT deals

127  with recently diverging paralogs by using neighborhood gene information. All the parameters

128  were set to default values except for the length ratio to discard shorter protein fragments

129  when a protein is split due to a frameshift or other mechanisms was set to 1.33 as

130  recommended by the authors. Orthologous clusters were stringently defined as all sequences

131  in a cluster having shared sequence identity ≥ 70 % and coverage ≥ 75 %. Plots and

132  calculations of pan-genome sizes, new genes discovered and pan-genome status were also

133  determined as described previously [21].

134  **Characterization of strains**.

135  *In silico* multilocus sequence typing of 198 strains was performed with the MLST 1.8 online

136  server [22]. The antimicrobial resistance genes in the sequenced isolates were identified by

137  BLASTp [23] searching with the databases of ARDB [24]. Genes conferring virulence

138  factors were identified using BLASTp with VFDB [25]. Given that many virulence factors

6

139     for *S. epidermidis* that are not contained in the VFDB, we used the orthologous proteins and

140     virulence factors from RP62a [26]  and ATCC1228 [27] to make up the missing information.

141     **Functional analysis**

142     All genes are BLASTed against all sequences in the database of KOBAS 2.0

143     (http://kobas.cbi.pku.edu.cn/) [28]. The cutoffs are BLASTp *E*-value <10-5 and BLAST

144     subject coverage > 70 %. We used the genes from same genome as the default background

145     distribution and considered only pathways for which there were at least two genes mapped.

146     Fisher's exact test was choosing to perform statistical test and Bonferroni correction was used

147     to reduce the high overall Type-I error with p.adjust from R package.

148     **Statistical analyses**

149     The differences in the prevalence of antimicrobial resistance genes and phenotypes among

150     isolates were analyzed by using two-tailed Fisher's exact test and Bonferroni correction was

151     also performed as mentioned above. All the statistical analyses were carried out using R

152     package (version: 3.3). A *P* value of < 0.05 was regarded as statistically significant.

153

154 **Results**

155 **Core pan-genome of *S. epidermidis***

156 Despite the intensive effort to characterize *S. epidermidis* and the sizable number of whole

157 genome comparisons in literature [29], more and more genome data is rapid accumulated and

158 could easily obtained from public database, such as NCBI. Using PanOCT, a total of 8,034

159 orthologous protein clusters were identified from a collection of all *S. epidiermidis* genomes

160 publicly available at the time of the analysis (Supplementary Table S1). PanOCT only

161 includes non-paralogs in clusters and uses conserved gene neighborhood to separate

162 duplicated genes. This means that insertion sequence elements that are in novel contexts will

163 often form singleton clusters even though they are identical in sequence to other IS elements

164 within or between genomes analyzed. When the "core" pan-genome is defined to be present

165 at all 198 genomes analyzed, there were 435 (5.4 %) core protein clusters and 2915 (36.3 %)

166 novel clusters (groups with a single member from a single genome) (Fig. 1a). To predict the

167 theoretical maximum pan-genome size (i.e., the total number of genes, including core, unique,

168 and dispensable genes) a pan-genome model was implemented using medians and an

169 exponential decay function (Fig. 1b). The maximum pan-genome size was estimated to be

170 12,554 ± 65 genes. To determine whether the *S. epidiermidis* pan-genome is open or closed,

171 the number of new genes identified (i.e., unique or strain-specific genes) for each genome

172 added was determined and fit to a power law function ($n = \kappa N^{-\alpha}$) as described previously

173 [21]. According to the result, we found the pan-genome of *S. epidiermidis* appeared to be

174 open ($\alpha = 0.226 \pm 0.002$; Fig. 1b). For each genome added, the number of new genes was

175 extrapolated by calculating tg(θ), which was determined to be 7.7 ± 0.4 (Fig. 1b).

176 The function of the genes within the variable genome was investigated by assigning all gene

177 clusters to clusters of orthologous groups (COGs) categories [30] and the results showed that

178 novel genes were most likely to be assigned to categories (Supplementary Table S2 and S3)

179    such as mobilome, ribosomal structure and biogenesis, carbohydrate transport and

180    metabolism, and nucleotide transport and metabolism, based on the result of Fisher's exact

181    test.

**Phylogenetic relationship of *S. epidermidis* isolates**

183    To estimate the genetic relationships among *S. epidermidis* strains, we compared all 198

184    genomes by using a single nucleotide polymorphism-based phylogeny. SNPs were identified

185    from the combined set of genome sequences by using kSNP. Nucleotide positions present in

186    at least 80 % of all genomes were used to build a Maximum-Likelihood phylogenetic tree

187    with RAxML following the tutorial. Strikingly, the 198 *S. epidermidis* isolates formed two

188    distinct groups (Fig. 3), called Cluster A (solid line) and B (dotted line). Most of Sequence

189    Type (ST) 2 nosocomial isolates were near identical at the nucleotide level for all core genes

190    (Supplementary Table S4). All of ST 2 strains in this study presented in Cluster A and had an

191    extremely short evolutionary distance from each other, indicating that these strains were

192    probably derived from a recent common ancestor. By contrast, Cluster B represents a lineage

193    with reduced virulence and all of ST 5 commensal strains presented in Cluster B and

194    clustered together. The rest of Cluster B had a much longer evolutionary distance from ST 5

195    strains. This clade may have more complex history of evolution and produce a various sub-

196    group.

**Antimicrobial resistance across *S. epidermidis***

198    Antimicrobial resistance (AMR) is very common among *S. epidermidis* isolates and often

199    limits treatment options [31]. Given the clinical importance of AMR in *S. epidermidis*, we

200    performed a genome-wide analysis of all known AMR genes within our genomic dataset.

201    According to the analysis of ARDB database, we found 28 different types of genes involved

202    in 31 antibiotics (Fig. 3). Nearly all isolates carry at least one type antibiotic resistance gene.

203    Among the genes involved in antimicrobial resistance, our data showed that there were two

204 genes, *nor*A and *bac*A, conserved in all strains. Based on the enrichment analysis of strains

205 from different sources, we found that strains from sources (skin, blood, environment and

206 plant) had significantly different antibiotic resistance profiles: isolates from blood (9

207 antibiotic resistance genes) and skin (8 antibiotic resistance genes) had significantly enriched

208 antibiotics (Supplementary Table S5), while isolates from environment had no significantly

209 enriched antibiotics. First-line antibiotic therapy for catheter-related bloodstream infections

210 was vancomycin. None of the isolates were resistant to the antibiotic at the genetic level,

211 regardless of isolation source.

**SCC*mec* and ACME in *S. epidermidis***

213 SCC*mec*, or *staphylococcal* cassette chromosome *mec*, is a mobile genetic element that

214 carries the central determinant for broad-spectrum beta-lactam resistance encoded by the

215 mecA gene a mobile genetic element of *Staphylococcus* bacterial species [10, 32]. According

216 to the completeness of genome in this study (only 7 complete genome sequences), we only

217 analyzed the genes from well-defined SCC*mec* genomic islands [33]. There were 58.6 %

218 (116/198) of *S. epidermidis* strains, in which complete *mec* gene complexes, *mec*A, and

219 *mec*R1 genes were detected (Supplementary Table S1). However, only 39.4 % (78/198) of

220 strains had *ccr* gene complex from type IV cassette, in which both *ccr*A and *ccr*B were

221 present. Similar to the previous results [29, 34], nearly all of the ST2 nosocomial isolates

222 (94.6 %, 70/74) had at least one copy of *mec*A from type IX cassette and *mec*R1 from Type

223 VIII or IV-G cassette. On the other hand, a high prevalence (98 %, 195/198) of ACME was

224 found in *S. epidermidis* strains in this study, of which 22.7 % (45/198) was type I and 75.8%

225 (150/198) was type II.

**Biofilm formation of *S. epidermidis***

227 Biofilm formation is the major of virulence factor of *S. epidermidis* strains, which will

228 contribute to the persistence of clinical infection. Here, we analyze some well-known genes

229    involved in biofilm formation such as adhesive molecules, including polysaccharide

230    intercellular adhesin (*ica*ABCD), proteinaceous factors (*bhp* and *aap*), teichoic acids,

231    extracellular DNA and so on (Supporting information Table S6). The polysaccharide

232    intercellular adhesion (*ica*ABCD) genes that encode biofilm-associated genes for poly-N-

233    acetylglucosamine synthesis were found in 60% of the commensal isolates, in agreement with

234    previous studies [34]. Especially, any of the *ica* genes was not found in some ST 2 strains

235    (Fig. 4). Gene *aap* was enriched in the blood (adjusted *P*-value < 0.01) compared to the

236    remaining isolates and therefore might be a potential biomarker for *S. epidermidis* infection.

237    We analyze the enrichment of all genes involved in virulence factors and found the *ica*ABCD

238    was significantly enriched despite the sources or sequence types.

**Human-Bacterium and Bacterium-bacterium interactions in *S. epidermidis***

240    *S. epidermidis* is the major colonization microorganisms in the human skin with complex

241    human-bacterium and bacterium-bacterium interactions. We analyzed the genes (Table 1 and

242    supplementary Table S5) involved in resistance against antimicrobial peptides that can inhibit

243    the growth of most skin microorganism including *S. epidermidis*. Some genes (e.g.

244    *cap*ABCD), which are significantly enriched in the blood and skin, were reported to assist the

245    strain to survive on the skin surface [7]. We also analyzed genes involved in bacterium-

246    bacterium interactions. We found that the genes involved in short-chain fatty acids

247    biosynthesis and extracellular proteases (e.g. Esp) had no difference despite the isolates.

248    Table 1 *S. epidermidis* resistance mechanisms that target AMPs.

| Resistance mechanism | Gene | Target AMPs | Functions | Enrichment |
|---|---|---|---|---|
| AMP sensing | *aps*SRX | Most cationic AMPs | 3-component sensor/regulator | - |
| | braSR/braDE/vraDE | Bacitracin, nisin | | - |
| Phosphatidylglycerol lysylation | *mpr*F | Most cationic AMPs | Lysylation of membrane phospholipids | - |
| Teichoic acids alanylation | *dlt*ABCD | Most cationic AMPs | Alanylation of teichoic acids | Blood / Skin (*dlt*D) |
| Exopolymers | *ica*ADBC | HBD3, LL-37, DCD-1 | Production of PNAG | Blood(*ica*B) |

| | | | exopolysaccharide; IcaB *N*-acetylglucosamine deacetylase introduces positive charge | |
|---|---|---|---|---|
| | *cap*ABCD | HBD3, LL-37, DCD-1 | | |
| Extracellular proteases | *sep*A | LL-37 | Degrades AMPs | - |
| | *esp* | LL-37a | | - |
| ABC transporters | *vra*FG | Vancomycin, polymyxin B, colistin | Putative AMP exporter | - |

249

250

## Discussion

251    *S. epidermidis* is a coagulase-negative and Gram-positive staphylococcus that is part of the

253    normal mucosa and skin microflora in humans and other mammals [2]. It is the second

254    leading cause of nosocomial infections [35]. Although it is a saprophyte, opportunistic

255    pathogen with plenty of antibiotic resistance and virulence factors [36], this natural skin

256    colonizer plays a critical role in balancing the epithelial microflora [1, 37]. As an innocuous

257    commensal microorganism, for a long time *S. epidermidis* has been seen as an avirulent

258    species. With the accumulation of genomic sequences, we can now further explore the

259    genetic mechanisms of environmental adaptability of *S. epidermidis*, the evolution process

260    during the outbreak, and the molecular biomarkers for clinic diagnosis [1, 38].

261    In our current pan-genome analysis, *S. epidermidis* had a relatively compact genome with a

262    size of about 2.5 Mb, and yet almost 20% of this genome was in flux, exchanging with a

263    large pool of various genes. These findings were similar to what had been reported by Conlan

264    and colleagues [29]. The significant number of genes involved in mobilome make horizontal

265    gene transfer easier between the *Staphylococcus* stains and lead to the increase of the "open"

266    pan-genome [39]. Besides, mobile genetic elements, such as SCC*mec*, ACME and plasmids,

267    make the genome structure more unpredictable [40]. High-resolution phylogenetic tree

268    constructed from genome-wide SNPs reveal important details not seen by traditional multi-

269    locus sequence typing (MLST) or single gene marker (16S rDNA). From the phylogenetic

270    tree, we found the ST2 isolates had an extremely short evolutionary distance from each other.

271    The genetic markers *mec*A and *ica*A, which are used to predict the antimicrobial resistance

272    and biofilm phenotypes, have been shown to be more common in hospital isolates than in

273    non-hospital isolates; however, these markers have much less power to distinguish infection

274    isolates from commensally available isolates that contaminate clinical specimens [41].

275    According to the enrichment analysis, we found it was impossible to distinguish the strains of

13

276  blood from that of skin, both of which had a similar lifestyle and genetic background.

277  However, it is possible to identify the strains from other habitats with biomarkers such as

278  *ica*ABCD and *cap*ABCD. Whole genome sequencing has been proved to be a more powerful

279  routine diagnostic tool than the traditional MLST or RT-PCR because it can rapidly identify

280  the infection source and antibiotic resistance in an affordable manner [42, 43]. As more

281  genetic data of *S. epidermidis* have been available and new machine learning algorithm is

282  developed [41], WGS may help to predict the infection isolation sources and antibiotic

283  resistance in a quicker and more accurate manner.

284  *S. epidermidis* has very complicated relationship with human and other bacteria.

285  Antimicrobial peptides (AMPs) play an important role in providing immunity to bacterial

286  colonization on human epithelia. Recent research has shown that staphylococci have multiple

287  systems to combat AMP activity, including AMP sensor that can regulate the expressions of

288  genes involved in AMP resistance depending on the presence of AMPs [7]. We analyzed the

289  distribution of gene involved in AMP resistance and found significant enrichment in blood

290  and skin and variable in different strains, which may be the consequence of coevolution of

291  human's immune system. On the other side, *S. epidermidis* strains also can inhibit the growth

292  of other bacterium to be dominant species on the skin surface. Serine protease Esp, which is

293  secreted by *S. epidermidis*, has been found to be able to inhibit the biofilm formation of *S.*

294  *aureus* and destroy pre-existing *S. aureus* biofilms [6]. Other mechanisms are also involved

295  in fighting against pathogens and maintaining homeostasis [44, 45]. On the other hand, *S.*

296  *epidermidis* was found to be a reservoir of antibiotic resistance, with its virulence

297  determinants shared with other more pathogenic species such as *S. aureus*, as demonstrated in

298  previous studies [29]. In particular, SCC*mec*, ACME elements conferring β-lactam resistance,

299  and other genes are transferred frequently between *Staphylococcus* strains, enabling rapid

300  evolution and adaptation against antibiotic selection pressure and provide additional

301 competitive advantage. For instance, type III of SCC*mec* carries a phenol soluble modulin

302 *psm-mec*, which may affect the virulence of *S. aureus* [40].

303 In conclusion, our current study provides information on the molecular characteristics of *S.*

304 *epidermidis* strains isolated from different environments from all over the world. From a

305 genomics perspective, the pan-genome analysis of the *S. epidermidis* reveals a high level of

306 diversity among the generic and species-specific genes and the potential supply routes for

307 enhanced versatility via inter- and intra-species horizontal gene transfer. Frequent horizontal

308 gene transfer enables the *Staphylococcus* to adapt to complex environments (e.g., high-level

309 antibiotic), and it may continue to be the dominant genus over the next few years. The

310 understanding of the mechanisms of gene transfer helps us to better prevent the emergence of

311 epidemic pan-drug resistant *S. epidermidis* strains.

312

313    **List of abbreviations**

314    SCC*mec*: Staphylococcal chromosome cassette *mec*

315    AMPs: Antimicrobial peptides

316    SNPs: Single-nucleotide polymorphisms

317    COGs: Clusters of orthologous groups

318    ST: Sequence type

319    AMR: Antimicrobial resistance

320

321

## **Acknowledgments**

328

329    References

330    1.    Otto M: *Staphylococcus epidermidis* - the "accidental" pathogen. *Nat Rev Microbiol*
331          2009, **7**(8):555-567.
332    2.    Oh J, Byrd AL, Park M, Program NCS, Kong HH, Segre JA: **Temporal Stability of the**
333          **Human Skin Microbiome**. *Cell* 2016, **165**(4):854-866.
334    3.    Schommer NN, Gallo RL: **Structure and function of the human skin microbiome**.
335          *Trends Microbiol* 2013, **21**(12):660-668.
336    4.    Jack RW, Tagg JR, Ray B: **Bacteriocins of gram-positive bacteria**. *Microbiol Rev* 1995,
337          **59**(2):171-200.
338    5.    Jetten AM, Vogels GD: **Mode of action of a Staphylococcus epidermidis bacteriocin**.
339          *Antimicrob Agents Chemother* 1972, **2**(6):456-463.
340    6.    Iwase T, Uehara Y, Shinji H, Tajima A, Seo H, Takada K, Agata T, Mizunoe Y:
341          *Staphylococcus epidermidis* **Esp inhibits** *Staphylococcus aureus* **biofilm formation**
342          **and nasal colonization**. *Nature* 2010, **465**(7296):346-349.
343    7.    Joo HS, Otto M: **Mechanisms of resistance to antimicrobial peptides in**
344          **staphylococci**. *Biochimica et biophysica acta* 2015, **1848**(11 Pt B):3055-3061.
345    8.    Foster TJ: **Antibiotic resistance in** *Staphylococcus aureus*. **Current status and future**
346          **prospects**. *FEMS Microbiol Rev* 2017.
347    9.    Diep BA, Gill SR, Chang RF, Phan TH, Chen JH, Davidson MG, Lin F, Lin J, Carleton HA,
348          Mongodin EF *et al*: **Complete genome sequence of USA300, an epidemic clone of**
349          **community-acquired meticillin-resistant** *Staphylococcus aureus*. *Lancet* 2006,
350          **367**(9512):731-739.
351    10.   McManus BA, Coleman DC, Deasy EC, Brennan GI, B OC, Monecke S, Ehricht R,
352          Leggett B, Leonard N, Shore AC: **Comparative Genotypes,** *Staphylococcal* **Cassette**
353          **Chromosome mec (SCC***mec***) Genes and Antimicrobial Resistance amongst**
354          *Staphylococcus epidermidis* **and** *Staphylococcus haemolyticus* **Isolates from**
355          **Infections in Humans and Companion Animals**. *PloS one* 2015, **10**(9):e0138079.
356    11.   Miragaia M, Thomas JC, Couto I, Enright MC, de Lencastre H: **Inferring a population**
357          **structure for Staphylococcus epidermidis from multilocus sequence typing data**. *J*
358          *Bacteriol* 2007, **189**(6):2540-2552.
359    12.   Bloemendaal AL, Brouwer EC, Fluit AC: **Methicillin resistance transfer from**
360          *Staphyloccccus epidermidis* **to methicillin-susceptible** *Staphylococcus aureus* **in a**
361          **patient during antibiotic therapy**. *PloS one* 2010, **5**(7):e11841.
362    13.   Yao Y, Sturdevant DE, Villaruz A, Xu L, Gao Q, Otto M: **Factors characterizing**
363          *Staphylococcus epidermidis* **invasiveness determined by comparative genomics**.
364          *Infection and immunity* 2005, **73**(3):1856-1860.
365    14.   Meric G, Miragaia M, de Been M, Yahara K, Pascoe B, Mageiros L, Mikhail J, Harris LG,
366          Wilkinson TS, Rolo J *et al*: **Ecological Overlap and Horizontal Gene Transfer in**
367          **Staphylococcus aureus and Staphylococcus epidermidis**. *Genome biology and*
368          *evolution* 2015, **7**(5):1313-1328.
369    15.   Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Lipman DJ, Ostell J, Sayers EW:
370          **GenBank**. *Nucleic Acids Res* 2013, **41**(Database issue):D36-42.
371    16.   Ugolotti E, Larghero P, Vanni I, Bandettini R, Tripodi G, Melioli G, Di Marco E, Raso A,
372          Biassoni R: **Whole-genome sequencing as standard practice for the analysis of**
373          **clonality in outbreaks of meticillin-resistant** *Staphylococcus aureus* **in a paediatric**
374          **setting**. *The Journal of hospital infection* 2016, **93**(4):375-381.

17. Barbier F, Lebeaux D, Hernandez D, Delannoy AS, Caro V, Francois P, Schrenzel J, Ruppe E, Gaillard K, Wolff M *et al*: **High prevalence of the arginine catabolic mobile element in carriage isolates of methicillin-resistant *Staphylococcus epidermidis*.** *J Antimicrob Chemother* 2011, **66**(1):29-36.

18. Gardner SN, Slezak T, Hall BG: **kSNP3.0: SNP detection and phylogenetic analysis of genomes without genome alignment or reference genome.** *Bioinformatics* 2015, **31**(17):2877-2878.

19. Stamatakis A: **RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.** *Bioinformatics* 2014, **30**(9):1312-1313.

20. Fouts DE, Brinkac L, Beck E, Inman J, Sutton G: **PanOCT: automated clustering of orthologs using conserved gene neighborhood for pan-genomic analysis of bacterial strains and closely related species.** *Nucleic Acids Res* 2012, **40**(22):e172.

21. Tettelin H, Masignani V, Cieslewicz MJ, Donati C, Medini D, Ward NL, Angiuoli SV, Crabtree J, Jones AL, Durkin AS *et al*: **Genome analysis of multiple pathogenic isolates of Streptococcus agalactiae: implications for the microbial "pan-genome".** *Proc Natl Acad Sci U S A* 2005, **102**(39):13950-13955.

22. Larsen MV, Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, Jelsbak L, Sicheritz-Ponten T, Ussery DW, Aarestrup FM *et al*: **Multilocus sequence typing of total-genome-sequenced bacteria.** *J Clin Microbiol* 2012, **50**(4):1355-1361.

23. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ: **Basic local alignment search tool.** *J Mol Biol* 1990, **215**(3):403-410.

24. Liu B, Pop M: **ARDB--Antibiotic Resistance Genes Database.** *Nucleic Acids Res* 2009, **37**(Database issue):D443-447.

25. Chen L, Xiong Z, Sun L, Yang J, Jin Q: **VFDB 2012 update: toward the genetic diversity and molecular evolution of bacterial virulence factors.** *Nucleic Acids Res* 2012, **40**(Database issue):D641-645.

26. Gill SR, Fouts DE, Archer GL, Mongodin EF, Deboy RT, Ravel J, Paulsen IT, Kolonay JF, Brinkac L, Beanan M *et al*: **Insights on evolution of virulence and resistance from the complete genome analysis of an early methicillin-resistant *Staphylococcus aureus* strain and a biofilm-producing methicillin-resistant *Staphylococcus epidermidis* strain.** *J Bacteriol* 2005, **187**(7):2426-2438.

27. Zhang YQ, Ren SX, Li HL, Wang YX, Fu G, Yang J, Qin ZQ, Miao YG, Wang WY, Chen RS *et al*: **Genome-based analysis of virulence genes in a non-biofilm-forming Staphylococcus epidermidis strain (ATCC 12228).** *Mol Microbiol* 2003, **49**(6):1577-1593.

28. Xie C, Mao X, Huang J, Ding Y, Wu J, Dong S, Kong L, Gao G, Li CY, Wei L: **KOBAS 2.0: a web server for annotation and identification of enriched pathways and diseases.** *Nucleic Acids Res* 2011, **39**(Web Server issue):W316-322.

29. Conlan S, Mijares LA, Program NCS, Becker J, Blakesley RW, Bouffard GG, Brooks S, Coleman H, Gupta J, Gurson N *et al*: ***Staphylococcus epidermidis* pan-genome sequence analysis reveals diversity of skin commensal and hospital infection-associated isolates.** *Genome Biol* 2012, **13**(7):R64.

30. Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin EV, Krylov DM, Mazumder R, Mekhedov SL, Nikolskaya AN *et al*: **The COG database: an updated version includes eukaryotes.** *BMC Bioinformatics* 2003, **4**:41.

31. Kleinschmidt S, Huygens F, Faoagali J, Rathnayake IU, Hafner LM: ***Staphylococcus epidermidis* as a cause of bacteremia.** *Future microbiology* 2015, **10**(11):1859-1879.

32. (IWG-SCC) IWGotCoSCCE: **Classification of *staphylococcal* cassette chromosome mec (SCCmec): guidelines for reporting novel SCCmec elements**. *Antimicrob Agents Chemother* 2009, **53**(12):4961-4967.

33. Kos VN, Desjardins CA, Griggs A, Cerqueira G, Van Tonder A, Holden MT, Godfrey P, Palmer KL, Bodi K, Mongodin EF *et al*: **Comparative genomics of vancomycin-resistant Staphylococcus aureus strains and their positions within the clade most commonly associated with Methicillin-resistant S. aureus hospital-acquired infection in the United States**. *mBio* 2012, **3**(3):e00112-00112.

34. Du X, Zhu Y, Song Y, Li T, Luo T, Sun G, Yang C, Cao C, Lu Y, Li M: **Molecular analysis of *Staphylococcus epidermidis* strains isolated from community and hospital environments in China**. *PloS one* 2013, **8**(5):e62742.

35. Ziebuhr W, Hennig S, Eckart M, Kranzler H, Batzilla C, Kozitskaya S: **Nosocomial infections by *Staphylococcus epidermidis*: how a commensal bacterium turns into a pathogen**. *Int J Antimicrob Agents* 2006, **28 Suppl 1**:S14-20.

36. Namvar AE, Bastarahang S, Abbasi N, Ghehi GS, Farhadbakhtiarian S, Arezi P, Hosseini M, Baravati SZ, Jokar Z, Chermahin SG: **Clinical characteristics of *Staphylococcus epidermidis*: a systematic review**. *GMS Hyg Infect Control* 2014, **9**(3):Doc23.

37. Otto M: ***Staphylococcus* colonization of the skin and antimicrobial peptides**. *Expert review of dermatology* 2010, **5**(2):183-195.

38. Didelot X, Bowden R, Wilson DJ, Peto TE, Crook DW: **Transforming clinical microbiology with bacterial genome sequencing**. *Nat Rev Genet* 2012, **13**(9):601-612.

39. Palmer KL, Kos VN, Gilmore MS: **Horizontal gene transfer and the genomics of Enterococcal antibiotic resistance**. *Curr Opin Microbiol* 2010, **13**(5):632-639.

40. Qin L, McCausland JW, Cheung GY, Otto M: **PSM-Mec-A Virulence Determinant that Connects Transcriptional Regulation, Virulence, and Antibiotic Resistance in *Staphylococci***. *Frontiers in microbiology* 2016, **7**:1293.

41. Tolo I, Thomas JC, Fischer RS, Brown EL, Gray BM, Robinson DA: **Do *Staphylococcus epidermidis* Genetic Clusters Predict Isolation Sources?** *J Clin Microbiol* 2016, **54**(7):1711-1719.

42. Rasko DA, Webster DR, Sahl JW, Bashir A, Boisen N, Scheutz F, Paxinos EE, Sebra R, Chin CS, Iliopoulos D *et al*: **Origins of the *E. coli* strain causing an outbreak of hemolytic-uremic syndrome in Germany**. *N Engl J Med* 2011, **365**(8):709-717.

43. Pankhurst LJ, Del Ojo Elias C, Votintseva AA, Walker TM, Cole K, Davies J, Fermont JM, Gascoyne-Binzi DM, Kohl TA, Kong C *et al*: **Rapid, comprehensive, and affordable mycobacterial diagnosis with whole-genome sequencing: a prospective study**. *Lancet Respir Med* 2016, **4**(1):49-58.

44. Wang Y, Kuo S, Shu M, Yu J, Huang S, Dai A, Two A, Gallo RL, Huang CM: ***Staphylococcus epidermidis* in the human skin microbiome mediates fermentation to inhibit the growth of *Propionibacterium acnes*: implications of probiotics in acne vulgaris**. *Appl Microbiol Biotechnol* 2014, **98**(1):411-424.

45. Otto M, Echner H, Voelter W, Gotz F: **Pheromone cross-inhibition between *Staphylococcus aureus* and *Staphylococcus epidermidis***. *Infection and immunity* 2001, **69**(3):1957-1960.

46. Grant JR, Arantes AS, Stothard P: **Comparing thousands of circular genomes using the CGView Comparison Tool**. *BMC Genomics* 2012, **13**:202.

469
470

471    Figure Legends

472    **Figure 1** Analysis of the *Staphylococcus epidermidis* pan-genome. (a). The distribution of

473    protein cluster sizes generated from the comparison of 198 *S. epidermidis* genomes using

474    PanOCT. (b). The pan-genome size (left) and the number of novel genes discovered with the

475    addition of each new genome (right) were estimated for all 198 genomes using a pan-genome

476    model based on the original Tettelin et al. model [21].

477    **Figure 2** Functional analysis of the pan-genome of *Staphylococcus epidermidis*. (a).

478    Distribution of core / dispensable / novel genes in the type strain RP62a. Starting from the

479    outermost ring the feature rings depict: (1) COG functional categories for forward strand

480    coding sequences; (2) Core (brown) / Dispensable (blue) genes for forward strand coding

481    sequences; (3) Forward strand sequence features; (4) Reverse strand sequence features; (5)

482    Core (brown) / Dispensable (blue) genes for reverse strand coding sequences; (6) COG

483    functional categories for reverse strand coding sequences. (7) GC content; (8) GC skew. The

484    colors of different COG functional categories were following the definition of Grant et al.

485    [46].

486    (b). Numbers of core, dispensable and novel genes for each COG category. COGs

487    significantly enriched (adjusted *P*-value < 0.05, Fisher exact test) in core, dispensable, or

488    novel genes are marked with red asterisk.

489    **Figure 3** Phylogenetic SNP tree of *Staphylococcus epidermidis* strains. A whole-genome

490    core SNP maximum likelihood tree was constructed for 198 genomes with kSNP and

491    RAxML. Heatmap on the right indicates copies of 28 genes involved in antibiotic resistance.

492    Legends on the bottom stand for copy number of resistant genes.

493    **Figure 4** Heatmap of virulence factors among the *Staphylococcus epidermidis* strains. The

494    dendrogram was generated using complete linkage clustering of copies of genes involved in

495    virulence factors. The red color stands for genes that exist in the genomes and the blue color

496    for missing ones. Legends on the right stand for colors of different host, isolates and

497    geographic information.

498    **Figure 5** *In silico* analysis of virulence factors of the *Staphylococcus epidermidis* strains. The

499    types of virulence factors were following the VFDBs database. Legends on the right stand for

500    colors of different host, isolates, and geographic information. Different colors stand for copy

501    number of each virulence factors.

502

503 Additional files

504 Table S1 Basic information of all strains analyzed used in this study

505 Table S2 Result of COG enrichment analysis across all strains

506 Table S3 Result of KEGG enrichment analysis across all strains
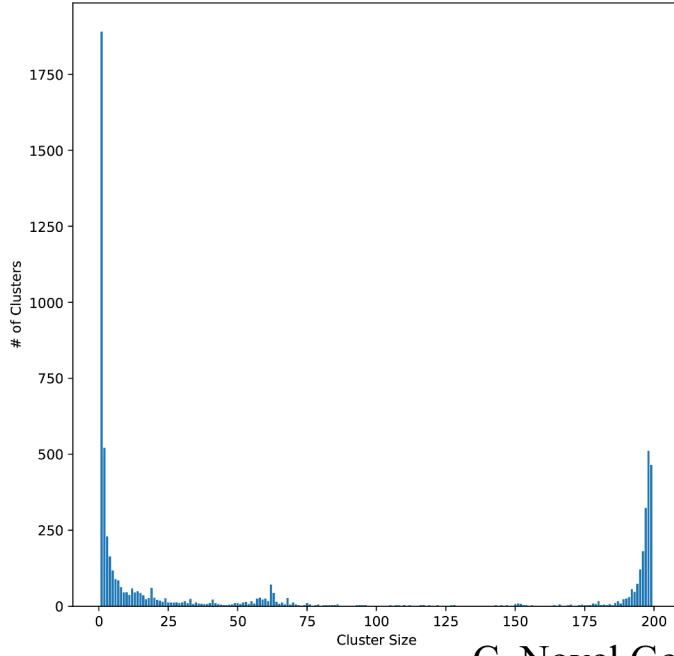
507 Table S4 Cluster of orthologous proteins produced by PanOCT

508 Table S5 Enrichment analysis about antibiotic resistances from different sources
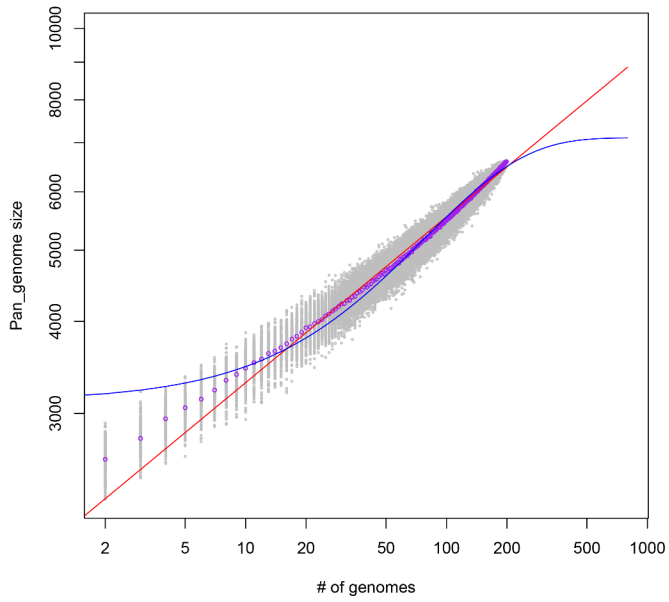
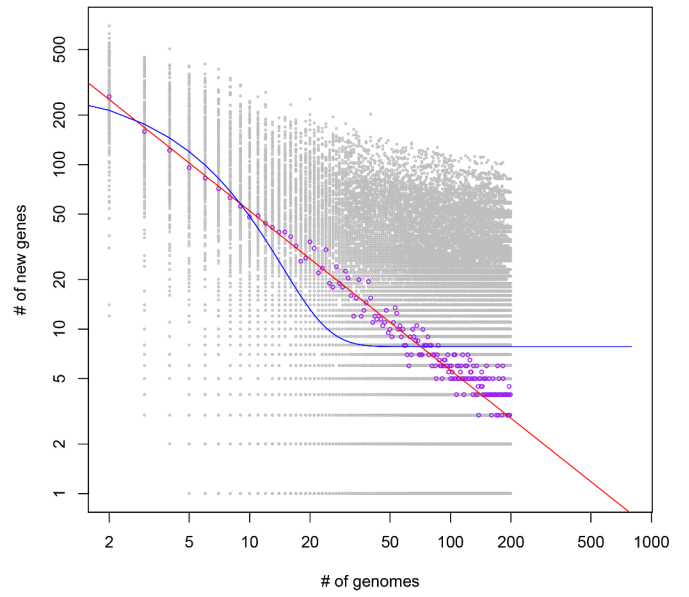509 Table S6 Enrichment analysis about genes related in biofilm formation
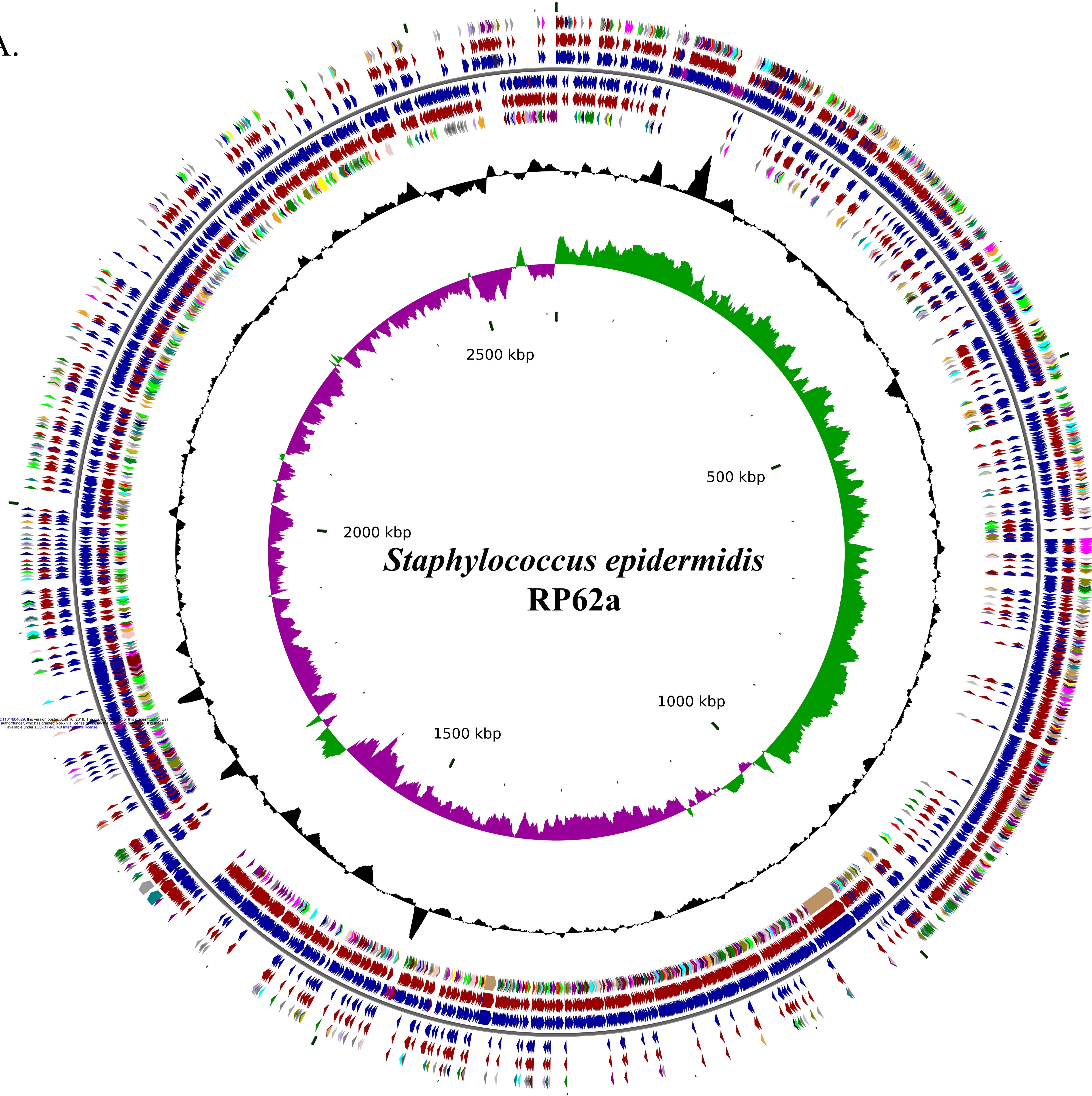
510

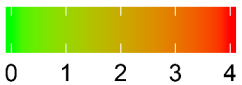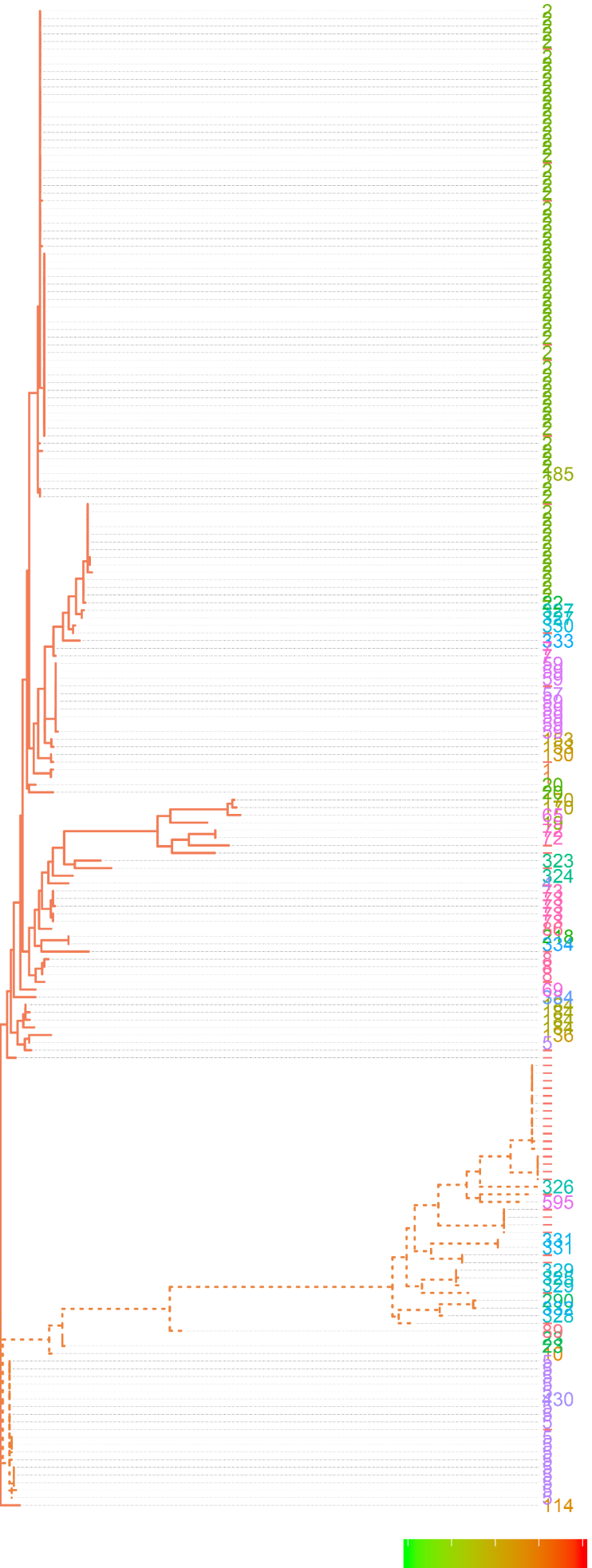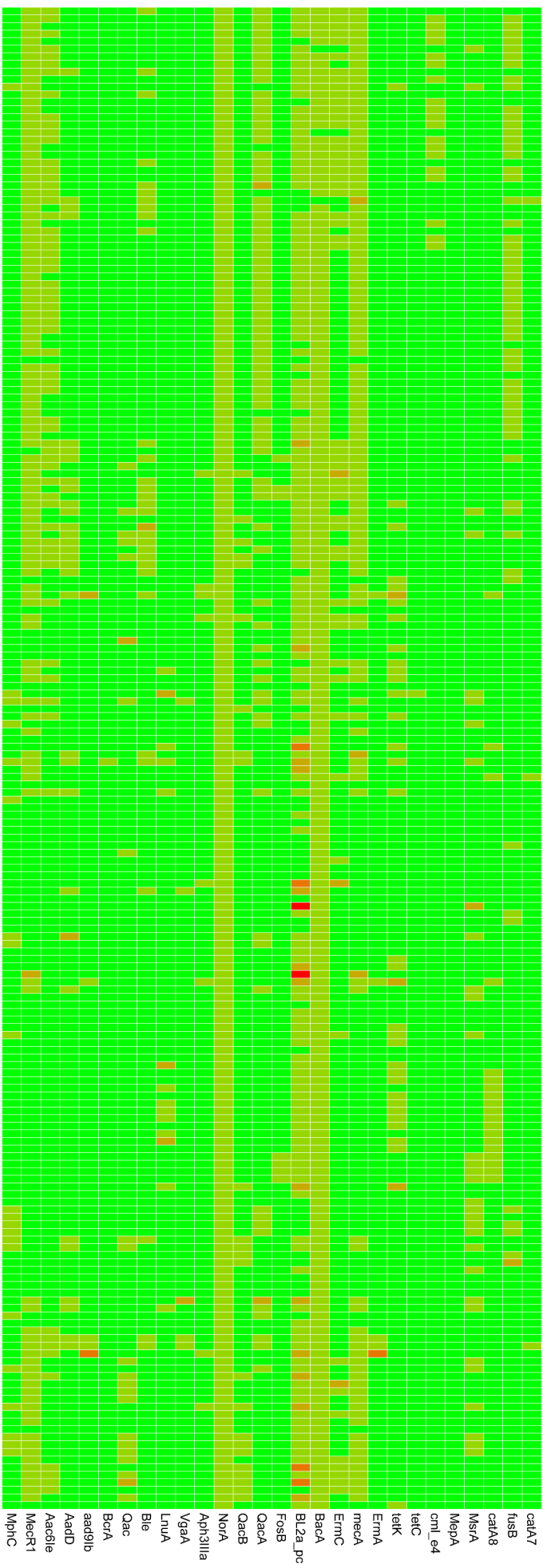# A. PanOCT Cluster Size Distribution



# B. Pan-Genome

# C. Novel Genes

A.

*Staphylococcus epidermidis*
RP62a

2500 kbp

2000 kbp

1500 kbp

1000 kbp

500 kbp

B.

Translation, ribosomal structure and biogenesis

Transcription

Signal transduction mechanisms

Secondary metabolites biosynthesis, transport and catabolism

Replication, recombination and repair

Posttranslational modification, protein turnover, chaperones

Nucleotide transport and metabolism

Mobilome: prophages, transposons    *

Lipid transport and metabolism

Intracellular trafficking, secretion, and vesicular transport

Inorganic ion transport and metabolism

General function prediction only

Function unknown

Energy production and conversion

Defense mechanisms

Coenzyme transport and metabolism

Cell wall/membrane/envelope biogenesis

Cell motility

Cell cycle control, cell division, chromosome partitioning

Carbohydrate transport and metabolism

Amino acid transport and metabolism

COG

Count

0    10000    20000    30000    40000
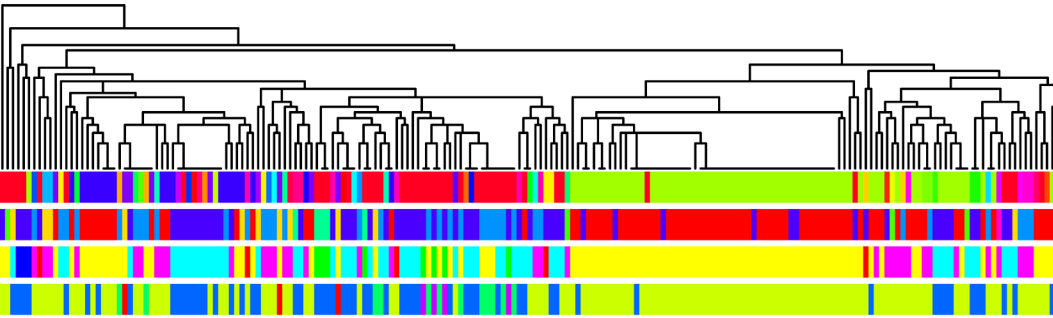
Type
Core
Dispensable
Novel

ST

Source

Location

Host

**Source**

Blood
Environ
Lung
Plants
Skin
Unknown
Urine

**Location**

Denmark
Germany
India
Other
Russia
USA

**Host**

Cow
Human
Mouse
Other
Rice

Capsule(capC)
Thermonuclease(nuc)
Polysaccharide_intercellular_adhesin(icaA)
Staphyloferrins
Antimicrobial_peptide_sensor(apsX)
Lipase(geh1)
Polysaccharide_intercellular_adhesin(icaC)
Polysaccharide_intercellular_adhesin(icaB)
Polysaccharide_intercellular_adhesin(icaD)
surface_protein_H(sesH)
surface_protein_G(sesG)
Serine-aspartate_repeat-containing_protein_H(sdrH)
Serine-aspartate_repeat-containing_protein_G(sdrG)
Phenol-soluble_modulins(hld)
d-alanylation_of_teichoic_acids(dltA)
Antimicrobial_peptide_sensor(apsR)
Serine-aspartate_repeat-containing_protein_F(sdrF)
eDNA(atlE)
surface_protein_E(sesE)
beta-hemolysin(hlb)
Teichoic_acids(tagA)
Esterase_1
Antimicrobial_peptide_sensor(apsS)
Biofilm-associated_protein_homolog(bhp)
d-alanylation_of_teichoic_acids(dltC)
d-alanylation_of_teichoic_acids(dltD)
Phenol-soluble_modulins_1
Phenol-soluble_modulins_2
Phenol-soluble_modulins_3
Phenol-soluble_modulins_4
Capsule(capB)
Accumulation-assocaited_protein(aap)
surface_protein_I(sesI)
eDNA(cidA)
Lipase(lipA)
Phenol-soluble_modulins_5
surface_protein_A(sesA)
Teichoic_acids(tagD)
Clp_Protease(clpX)
protease(Serine)
Clp_Protease(clpB)
Teichoic_acids(tagX)
Teichoic_acids(tagH)
vraFG(vraG)
Lipase(geh2)
Clp_Protease(clpC)
Autolysin/adhesin(aae)
vraFG(vraF)
Capsule(capD)
Teichoic_acids(tagB)
Serine_V8_protease(sspA)
Multiple_peptide_resistance_factor(mprF)
Lipase(lip)
Capsule(capA)
Lipase(geh)
Hemolysin
Extracellular_matrix_binding_protein(ebh)
d-alanylation_of_teichoic_acids(dltB)
Phenol-soluble_modulins_6
Esterase_2
Hemolysin_III
Zinc_metalloprotease
surface_protein_C(sesC)
Cysteine_protease(sspB)
Cysteine_protease(sspC)
Clp_Protease(clpP)
Teichoic_acids(tagG)
Nuclease
Iron_transporter(sitA)
Elastase(sepA)
Iron_transporter(sitC)
Iron_transporter(sitB)

Source
- Blood
- Environment
- Lung
- Plants
- Skin
- Unknown
- Urine

Locatioin
- Denmark
- Germany
- India
- Other
- Russia
- USA

Host
- Cow
- Human
- Mouse
- Other
- Rice

Source
Locatioin
Host

Intercellular adhesion proteins

X delta hemolysin

ClpP

Aureolysin

Type VII secretion system

Lipase

SDr

Enterotoxin like L

Enterotoxin C

Staphopain

ClpC