1 # Individual level predictions of *Staphylococcus aureus* bacteraemia-associated
2 mortality.
3

4 Mario Recker[1*], Maisem Laabei[2*], Michelle S. Toleman[3], Sandra Reuter[3], Beth Blane[3], M.

5 Estee Török[3], Sion Bayliss[2], Sharon J. Peacock[3,4] and Ruth C. Massey[2ø].

6

7 1: Centre for Mathematics & the Environment, University of Exeter, Penryn Campus, Penryn

8 TR10 9EZ

9 2: Dept. of Biology and Biochemistry, University of Bath, Bath, BA2 7AY, UK

10 3: Department of Medicine, University of Cambridge, Cambridge, UK

11 4: London School of Hygiene and Tropical Medicine, London, WC1E 7HT, UK

12 * contributed equally to this work.

13 ø Corresponding author: r.c.massey@bath.ac.uk

14

15
16
17
18
19

20

21

22

23

24

25

26

# Abstract

The bacterium *Staphylococcus aureus* is a major human pathogen, where the emergence of antibiotic resistance is a global public-health concern. Host factors such as age and the presence of co-morbidities have been implicated in a worse outcome for patient. However, this is complicated by the highly complex and multi-faceted nature of bacterial virulence, which has so far prevented a robust mapping between genotype, phenotype and infection outcome. To investigate the role of bacterial and host factors in contributing to *S. aureus* bacteraemia-associated mortality we sequenced a collection of clinical isolates (of the MLST clonal complex CC22) from patients with bloodstream infections and quantified specific virulence phenotypes. A genome-wide association scan identified several novel virulence-affecting loci, which we validated using a functional genomics approach. Analysing the data comprising bacterial genotype and phenotype as well as clinical meta-data within a machine-learning framework revealed that mortality associated with CC22 bacteraemia is not only influenced by the interactions between host and bacterial factors but can also be predicted at the individual patient-level to a high degree of accuracy. This study clearly demonstrates the potential of using a combined genomics and data analytic approach to enhance our understanding of bacterial pathogenesis. Considering both host and microbial factors, including whole genome sequence and phenotype data, within a predictive framework could thus pave the way towards personalised medicine and infectious disease management.

# 1 Introduction

2 *Staphylococcus aureus* bacteraemia (SAB) is a significant global health problem[1] and is

3 exacerbated by the emergence and widespread circulation of drug resistant strains, such as

4 methicillin-resistant *S. aureus*, or MRSA[2]. Mandatory surveillance of SAB has now been

5 implemented in several countries, with many reporting a decline in the incidence of

6 methicillin-resistant SAB (MR-SAB)[3-5]. However, in the UK the incidence of methicillin-

7 susceptible SAB (MS-SAB) has been increasing year on year, where there has been a 15.4%

8 increase in cases since reporting became mandatory in 2011/2012[3]. Furthermore, the 30

9 day (all-cause) mortality rate for SAB has not significantly changed over the last two decades

10 and appears to have plateaued at approximately 20%[6]. This strongly suggests that existing

11 infection control and treatment options are insufficient to tackle this important health

12 problem and that a better understanding of the factors that contribute to bacteraemia-

13 associated morbidity is crucially needed.

14

15 To date, many host risk factors have been identified for both the occurrence of and

16 treatment failure following SAB[6]. However, the contribution of the bacterium is only

17 partially understood and is largely informed by experimental animal studies. These model

18 systems come with their own set of limitations, and many observations from these contrast

19 with those from human studies. For example, whereas cytolytic toxins have previously been

20 shown to enhance disease severity in animal models of SAB[7,8], isolates from invasive

21 diseases in humans, such as bacteremia and pneumonia, were recently found to be

22 significantly less toxic than those isolated from skin and soft tissue infections or even those

23 of healthy volunteers[9-12]. This raises the question as to whether animal models are adequate

1    to study bacterial virulence in human SAB infections, or whether there is an important

2    distinction between the role of toxicity in causing bacteremia and its pathogenic effect once

3    bacteremia has been established. Either way, human-based approaches are essential to

4    close this gap in our knowledge.

5

6    Another limitation to our understanding of the pathogenesis of SAB is that most studies

7    focus on only on a single or small number of factors in isolation, host or bacterial.  For

8    example, several studies have found increased mortality rates associated with MR-SAB

9    compared to MS-SAB[13,14]. However, patients with co-morbidities are more likely to develop

10   an MR-SAB due to their impaired health and longer time spent in healthcare facilities when

11   compared to those without comorbidities. When subsequent studies considered both

12   factors together, no difference in mortality was associated with the methicillin resistance

13   status of the infecting bacterium[15]. This illustrates the importance of a more inclusive

14   approach that considers all of the potential host and bacterial factors and examines how

15   their interactions influence the outcome for the patients.

16

17   We have previously demonstrated how genotype-phenotype mapping in *S. aureus* has the

18   potential to provide sufficient information to enable predictions of the level of virulence

19   expressed by the infecting microorganism[16]. Here, by expanding this whole-genome

20   approach to a set of fully sequenced isolates from bacteraemic patients and analysing the

21   combined set of genotype, phenotype and clinical meta-data within a machine learning

22   predictive modelling framework we show how host and bacterial factors interact to

4

1    determine severe infection outcome of *S. aureus* bacteraemia. These findings pave the way

2    towards individual patient-level predictions and personalised treatment strategies.

3

## Material and Methods.

5    **Strain and clinical metadata collection**. All isolates were collected from adults admitted to a

6    single hospital with their first episode of SAB between 2006 and 2012, and were stored in

7    glycerol at -80°C. Samples were collected during an observational cohort study of adults

8    with SAB at Addenbrooke's Hospital, Cambridge, UK between 2006 and 2012. Written

9    informed consent was not required as the study was conducted as part of a service

10   evaluation of the management of SAB. Ethical approval was obtained from the University of

11   Cambridge Human Biology Research Ethics Committee (reference HBREC.2013.05) and the

12   Cambridge University Hospitals NHS Foundation Trust Research and Development

13   Department (reference A092869). Study definitions have been defined previously[17] and

14   were used to determine the focus of the bacteremia, classify the bacteremia as community-

15   acquired, hospital-acquired or healthcare-associated, and to report outcomes, including

16   death at 30 days. The presence of comorbidities was assessed using the Charlson

17   comorbidity index (CCI) and were dichotomized into scores of <3 or ≥3 [18,19], as detailed in

18   Supplementary Table 1.

19   **Genome Sequencing**

20   Bacterial DNA extraction was carried out on a QIAxtractor (Qiagen), and library preparation

21   was performed as previously described[20]. Index-tagged libraries were created, and 96

22   separated libraries were sequenced in each of eight channels using the Illumina HiSeq

23   platform (Illumina) to generate 100-bp paired-end reads at the Wellcome Trust Sanger

1 Institute, UK. Paired-end reads for these isolates were mapped to the ST22/EMRSA15

2 reference strain, HO 5096 0412[21] and SNPs were identified as described previously[21]. The

3 accession number for the sequence data for each of these isolates is listed in Supplementary

4 Table 1.

5 **Cytotoxicity Assay**

6 Overnight *S. aureus* cultures were diluted 1:1000 into fresh tryptone soya broth and

7 incubated for 18 h at 37°C with shaking at 180 rpm in 30 mL glass tubes. *S. aureus*

8 supernatants were harvested from 18 h cultures by centrifugation at 14,600 rpm for 10 min.

9 The THP-1 human monocyte-macrophage cell line (ATCC#TIB-202) was routinely grown in

10 suspension in 30 mL of of RPMI-1640, supplemented with 10% heat-inactivated fetal bovine

11 serum (FBS), 1 μM L-glutamine, 200 units/mL penicillin, and 0.1 mg/mL streptomycin at

12 37°C in a humidified incubator with 5% $CO_2$. Cells were routinely viewed microscopically

13 every 48–60 h and harvested by centrifugation at 1,000 rpm for 10 min at room

14 temperature and re-suspended to a final density of 1–1.2 x $10^6$ cells/mL in tissue-grade

15 phosphate buffered saline. This procedure typically yielded >95% viability of cells as

16 determined by trypan blue exclusion and easyCyte flow cytometry. To evaluate *S. aureus*

17 toxicity, we diluted the supernatant to 30% of the original volume in TSB and incubated 20

18 μL of washed THP-1 cells with 20 μL diluted bacterial supernatant for 12 min at 37°C under

19 static conditions. Cell death was quantified by easyCyte (Millipore) flow cytometry using the

20 Guava Viability stain (Millipore) according to manufacturer's instructions, with the toxicity

21 of each isolate quantified in triplicate and the mean of this presented (listed in

22 Supplementary Table 1).

23 **Biofilm assay**

1    Biofilm formation was quantified using a 1:40 dilution from overnight cultures into 100 µL of

2    fresh tryptic soy broth supplemented with 0.5% sterile filtered glucose (TSBG) in 96-well

3    polystyrene plate (Costar). Perimeter wells of the 96-well plate were filled with sterile $H_2O$

4    and plates were placed in a separate plastic container inside a 37°C incubator and grown for

5    24 h under static conditions. For the transposon mutants, erythromycin (5 µg/mL) was

6    added to the growth medium. Semi-quantitative measurements of biofilm formation on 96-

7    well polystyrene plates was determined based on the method of Ziebuhr et al[22]. Following

8    24-h growth, plates were washed vigorously five times in PBS, dried and stained with 150 µL

9    of 1% crystal violet for 30 min at room temperature. Following five washes of PBS, wells

10   were re-suspended in 200 µL of 7% acetic acid, and optical density at 595 nm was recorded

11   using a Fluorimeter plate reader (BMG Labtech). To control for day to day variability, for the

12   clinical isolates a control strain (E-MRSA15) was included on each plate in triplicate, and

13   absorbance values were normalised against this (listed in Additional Table 1). For the

14   transposon mutants, as JE2 is the wild type strain this was used as the control strain, and

15   the effect of mutating the loci made relative to this. For this experiment the assays were

16   performed in triplicate on each plate and repeated four times.

17   **Genome wide association study (GWAS)**

18   The genome of the reference strain HO 5096 0412 was split into 2095 variable loci,

19   corresponding to coding region and intergenic regions, containing SNPs relative to the

20   reference genome. Annotated intergenic elements such as miscellaneous RNAs were

21   considered separate loci. Synonymous SNPs in coding regions and SNPs in known mobile

22   genetic elements and repeat regions were not considered. The resulting loci were named

23   allele_X where X refers to the position of the SNP at the 5' end of that block, relative to the

1    origin of replication. Each of these alleles in each isolate was scored as 1 if it differed from

2    the reference and 0 if it didn't. These allele scores for each isolate were used as the

3    genotypic information for the following analysis. Significant associations between bacterial

4    genotype and either phenotype (toxicity and biofilm formation) were identified by fitting an

5    analysis of variance model (ANOVA) in R[23] and using a minor allele frequency cut-off of 5%.

6    The reported *P* values are not corrected for multiple testing; the Bonferroni statistical

7    significance threshold is instead provided in fig. 2.

8    **Predictive modelling.**

9    We employed a *Random Forests*[24] machine learning approach, using the *randomForest*

10   package in R[25], to identify predictive signatures of host mortality based on the genotype,

11   phenotype and clinical meta-data. To assess the models' predictive accuracies we employed

12   two different measures: (i) the receiver operating characteristic (ROC) curve, which is

13   generated by plotting the true positive rate against the false positive rate (i.e. the observed

14   incidence against the false predicted incidence) at various threshold settings and where the

15   area under the curve (AUC) is a measure of predictive accuracy, with an AUC=1 equating to

16   zero error and an AUC=0.5 equating to random guessing; and (ii) by means of a confusion

17   matrix, which contrasts the instances of the predicted classes (*alive* or *death*) against the

18   actually observed classes. The misclassification-rates reported here are based on the so-

19   called out-of-bag errors[25], which are derived by iteratively testing the models' performances

20   against subsets of data left out during the fitting processes and thus provide a measure of

21   how well the models would fare against unknown data.

22

23

# Results

To elucidate the role of bacterial factors in human SAB we analysed a collection of 135 *S. aureus* isolates sampled from patients with bloodstream infections admitted to a single hospital between 2006 and 2012. All isolates belonged to the multi-locus defined clonal complex 22 (CC22) and contained both methicillin-resistant (MRSA) and methicillin-susceptible (MSSA) isolates. The 30-day all-cause mortality rate was 24.1% and the clinical data are summarized in Supplementary Table 1. Each isolate was whole-genome sequenced, and quantitatively phenotyped with respect to cytolytic activity (the ability to lyse a monocyte cell line, THP-1) and biofilm formation, both of which are major virulence determinants implicated in *S. aureus* disease[1,2]. Cytolytic toxins enable the evasion of cellular aspects of host immunity, release nutrient from host cells and are responsible for much of the purulent tissue damage associated of *S. aureus* infections[1,2]. Biofilm formation enables *S. aureus* to colonise foreign material and medical devices, protects the bacteria from many aspects of host immunity and renders some antibiotics less effective[1,2]. Despite the close genetic and geographic relationship between these isolates, toxicity and biofilm formation varied widely across the isolates; no association was found between methicillin susceptibility and either virulence phenotype (fig. 1).

We performed a genome-wide association study (GWAS) to identify loci that were associated with the toxicity of the isolates and their ability to form biofilm, where associations were tested at an uncorrected ($P<0.05$) and a Bonferroni corrected ($P<4.6 \times 10^{-5}$) significance threshold (fig. 2a). For toxicity, no loci reached significance when using the Bonferroni correction for multiple comparisons, although mutations in the Agr toxicity-regulating locus and a putative membrane protein (SAEMRSA15_10750) were marginally

1    associated (0.05 > $P$ > 4.6x10$^{-5}$). For biofilm, only one locus reached statistical significance,

2    *pbp2*, a penicillin binding protein. This gene is believed to be essential to the bacteria, as no

3    transposon mutants were available to functionally verify its effect on biofilm. However,

4    members of this family of proteins have been shown previously to affect biofilm

5    formation[26]. A further 37 loci showed marginally significant associations with biofilm

6    (Supplementary Table 2), of which we chose a subset of 12 for functional validation using

7    transposon mutants[27] (Fig. 2b). Out of these, transposon mutagenesis of six loci showed a

8    significant effect on biofilm formation compared to the wild type: a putative helicase

9    (SAEMRSA15_23880, *Tn* mutant NE513), the quinolone efflux protein NorA

10    (SAEMRSA15_16820, Tn mutant NE1034), a putative thiamine pyrophosphate

11    enzyme/indole-3-pyruvate decarboxylase (SAEMRSA15_01530, *Tn* mutant NE1149), and a

12    putative peptidase (SAEMRSA15_04760, *Tn* mutant NE1455), a putative inosine-uridine

13    preferring nucleoside hydrolase (SAEMRSA15_02020, *Tn* mutant NE1637, and a hypothetical

14    protein (SAEMRSA15_11850, *Tn* mutant NE318) (fig. 2b).

15

16    Having access to the genetic and phenotypic data for each isolate as well as clinical data

17    (Supplementary Table 1) for 92 of the 135 patients, we sought to determine the factors

18    most predictive of severe infection outcome, namely host mortality at day 30. Due to the

19    high dimensionality and high number of possible interactions between the various host and

20    bacterial factors we employed a *random forests* machine learning approach[24,25]. We

21    compared four scenarios, where the model was fitted against (i) the bacterial genotype data

22    only, (ii) the bacterial phenotype data only, (iii) the clinical data only, and (iv) the entire

23    dataset. Where genotype data were incorporated we first performed a feature selection

24    process, in which the model was fitted against the entire dataset, and a smaller subset of

1   variables selected for further analysis based on their predictive importance, i.e. their

2   contribution to the model's performance in distinguishing between the two classes. This

3   reduced the total number of variables used for fitting down to just 20 (from over 2000

4   initially).

5

6   As shown by receiver operator characteristic (ROC) curves in figure 3a, data comprising

7   either just bacterial factors (SNP or phenotype) or just clinical factors were relatively poor

8   predictors of mortality. The most accurate model was one that included a combination of all

9   factors, suggesting that these may interact in determining infection outcome. In this case,

10  the model yielded a remarkable individual patient-level predictive accuracy of nearly 80%

11  (based on out-of-bag error rates (which is a measure of prediction error)). This is shown by

12  the illustrated confusion matrix in fig. 3b, where the dark-blue diagonal sectors represent

13  the true positive and true negative rates and the light-blue off-diagonal entries represent

14  the false positive and false negative rates.

15

16  *Random forest* can also be used to rank factors based on their relative importance to the

17  predictive power of the model. As illustrated in fig. 3c toxicity and biofilm were the most

18  important features, with patient age and co-morbidities (CCI-class) also playing important

19  roles. Our model also identified a number of other bacterial genetic loci as contributing to

20  disease outcome (listed in Supplementary Tables 3). Of particular note is the *capA* gene

21  (indicated as allele_142543 in fig. 3c), which encodes a capsule biosynthesis enzyme. While

22  the role of capsule in protecting bacteria against several aspects of the host immune system

23  during experimental infections has been established[28,29], several clinically important clones

24  have evolved to be capsule-negative, due to polymorphisms in the capsule biosynthesis

1  genes[30]. The protective activity of capsule is therefore not critical to the ability of *S. aureus*

2  to cause disease in humans. Although we have not functionally verified the effect of the

3  polymorphism found here on capsule expression, it suggests that this bacterial phenotype

4  can also vary within a clone and potentially affect infection outcome.

5

6  Based on the ranked importance scores of the different factors used in this model (fig. 3c),

7  both toxicity and biofilm formation are highly influential in determining SAB-associated

8  mortality, the former positively and the latter negatively. To further analyse these and

9  illustrate their influence alongside the other important factors we used the *random forest*

10  approach again to derive two-way partial-dependence plots. As shown in figure 4, the

11  combination of age with high levels of toxicity and/or low levels of biofilm was associated

12  with a significant increase in the risk of patient death. Equally, the comorbidity index

13  (CCI)[18,19] with any combination of old age, high toxicity, and low levels of biofilm formation

14  considerably increased the probability of severe infection outcomes.

15

16  **Discussion**

17  In this study we demonstrate how combining genotypic, phenotypic and clinical metadata

18  within a mathematical framework can be used to accurately predict the mortality of

19  individual patients following SAB. This work also highlighting the critical role that specific

20  bacterial phenotypes play in affecting the outcome from SAB, providing clarity to an

21  apparent contradiction between experimental animal and human studies. Animal models

22  have demonstrated a clear role for toxins in the severity and disease outcome from SAB[7,8],

23  however work by us and several other groups have found a negative correlation between

24  toxicity and invasive diseases such as SAB in humans[9-12]. Here our data appears to conflict

1  with our previous work and support the findings of the animal-based studies, which we

2  believe can be reconciled by proposing a sequential pathway of events that distinguishes

3  between the establishment of SAB (i.e. gaining access to the bloodstream) and the

4  subsequent processes associated with a fatal outcome once SAB has become established.

5  Toxin production appears to play disparate roles in these two stages of the disease, being

6  selected against for the initial process of gaining access to the blood stream, but once in the

7  bloodstream it is a major factor leading to host mortality, likely due to the damage and

8  inflammation it causes there. A likely explanation for why animal models of SAB have only

9  confirmed the latter process (disease severity) is that the natural infectious process of

10  gaining access to the bloodstream is bypassed by injecting high doses (typically $10^7$-$10^8$

11  colony forming units) directly into the animal's bloodstream.

12

13  From the host perspective, the importance of patient age and co-morbidities in contributing

14  to mortality is consistent with previous findings. However, our modeling approach also

15  demonstrates that these host features can critically interact with bacterial factors and,

16  importantly, allow a prediction of patient outcome with a high degree of accuracy. It is likely

17  that other factors not considered here may influence a patient's risk of death following SAB.

18  Apart from patient care, host genotype and other bacterial phenotypes could have a

19  significant effect on infection outcome. A larger more detailed dataset may enable us to

20  fully identify these factors and unravel their interactions to predict mortality with even

21  higher accuracy. However, given the known epidemiology (with older patients with co-

22  morbidities being the most susceptible) and the complexity of this type of disease, it is

23  unlikely that we will achieve 100% predictability, demonstrating the significance of the near

24  80% accuracy reported here.

1

2   With the growing global problem of antimicrobial resistance, alternative intervention and

3   control strategies are needed. These include the development of vaccines and identification

4   of drugs that attenuate the virulence of pathogens. However, without a full understanding

5   of how the bacterial targets for these strategies are involved in causing disease in humans,

6   there is a significant risk of investing in and pursuing unsuccessful lines of therapeutic

7   development. Our findings here, for example, suggest that cytolytic toxins, components of

8   biofilm and possibly capsule are unlikely to be good targets, as they play disparate roles in

9   different stages of disease and their expression is highly variable, even within closely related

10  bacterial clones. Of greater importance, however, it that this work has the potential to make

11  a significant contribution to infectious disease diagnosis and management. With the move

12  towards the introduction of microbial genome sequencing into routine diagnostic settings,

13  the ability to use such information to inform clinicians on the likely outcome of an infection

14  for individual patients would allow them to tailor treatment to that individual, an important

15  step towards the implementation of personalised medicine approaches to infectious disease

16  management.

17

18  **References.**

19  1.  Lowy FD. *Staphylococcus aureus* infections. *N. Engl. J. Med.* 1998;339:520-32.

20  2.  Gordon, R.J. & Lowy FD. Pathogenesis of methicillin-resistant Staphylococcus aureus

21      infection. Clin Infect Dis. 2008;46:S350-9.

22  3.  https://www.gov.uk/government/uploads/system/uploads/attachment_data/file/53563

23      5/AEC_final.pdf

4.  Okon KO, Shittu AO, Kudi AA, Umar H, Becker K, Schaumburg F. Population dynamics of Staphylococcus aureus from Northeastern Nigeria in 2007 and 2012. Epidemiol Infect. 2014;142:1737-40.

5.  Walter J, Haller S, Blank HP, Eckmanns T, Abu Sin M, Hermes J. Incidence of invasive meticillin-resistant Staphylococcus aureus infections in Germany, 2010 to 2014. Euro Surveill. 2015;20(46).

6.  van Hal SJ, Jenson SO, Vaska VL, Espedido BA, Pateron DL, Gosbell IB. Predictors of mortality in Staphylococcus aureus Bacteremia. Clin Microbiol Rev. 2012;25:362-86.

7.  Jenkins A, Diep BA, Mai TT, Vo NH, Warrener P, Suzich J, et al. Differential expression and roles of Staphylococcus aureus virulence determinants during colonization and disease. MBio. 2015;6:e02272-14.

8.  Crémieux AC, Saleh-Mghir A, Danel C, Couzon F, Dumitrescu O, Lilin T, et al. α-Hemolysin, not Panton-Valentine leukocidin, impacts rabbit mortality from severe sepsis with methicillin-resistant Staphylococcus aureus osteomyelitis. J Infect Dis. 2014;209: 1773-80.

9.  Sharma-Kuinkel BK, Wu Y, Tabor DE, Mok H, Sellman BR, Jenkins A, et al. Characterization of alpha-toxin hla gene variants, alpha-toxin expression levels, and levels of antibody to alpha-toxin in hemodialysis and postsurgical patients with Staphylococcus aureus bacteremia. J Clin Microbiol. 2015;53:227-36.

10. Laabei M, Uhlemann AC, Lowy FD, Austin ED, Yokoyama M, Ouadi K, et al. Evolutionary Trade-Offs Underlie the Multi-faceted Virulence of Staphylococcus aureus. PLoS Biol. 2015;13:e1002229.
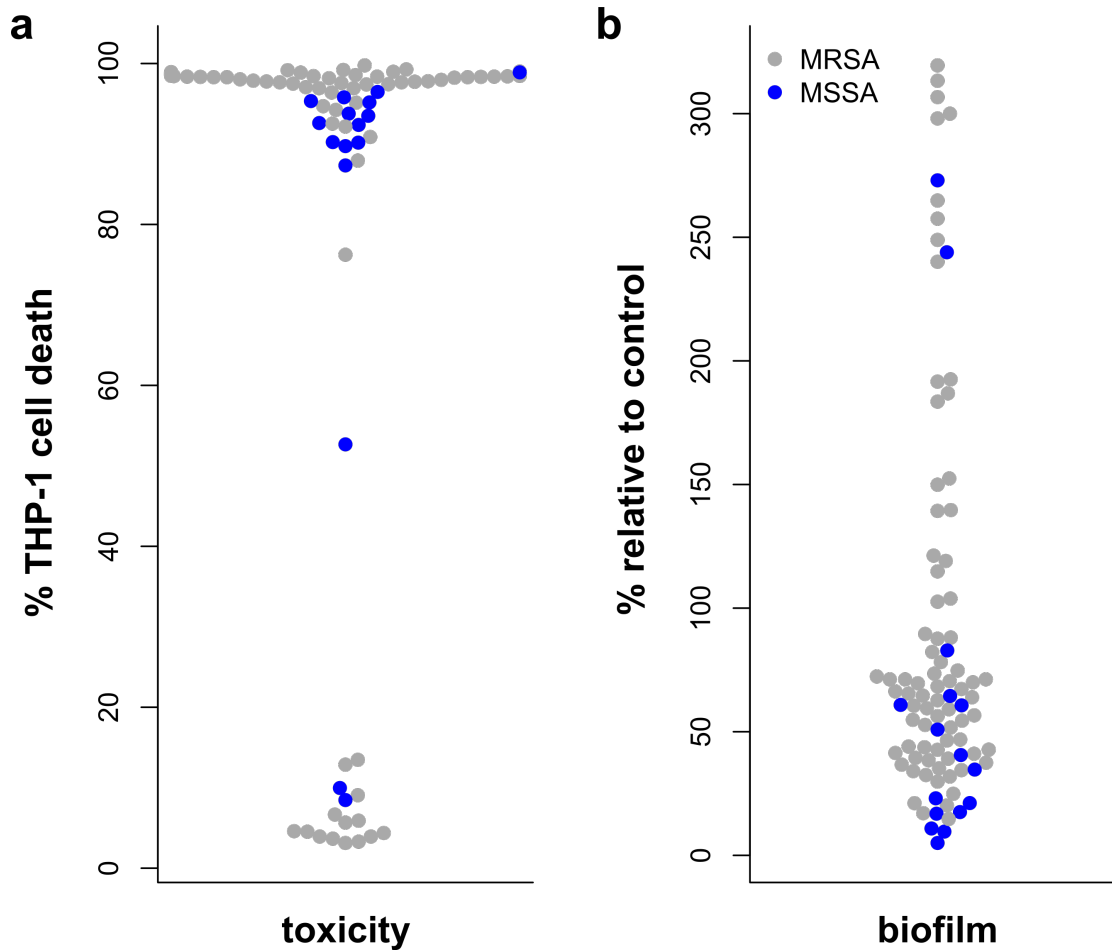
1    11. Rose HR, Hilzman RS, Altman DR, Smyth DS, Wasserman GA, Kafer JM, et al. Cytotoxic

2        Virulence Predicts Mortality in Nosocomial Pneumonia Due to Methicillin-Resistant

3        Staphylococcus aureus. J Infect Dis. 2015;211:1862-74.

4    12. Das S, Lindemann C, Young BC, Muller J, Osterreich B, Ternette N, et al. Natural

5        mutations in a Staphylococcus aureus virulence regulator attenuate cytotoxicity but

6        permit bacteremia and abscess formation. Proc Natl Acad Sci USA. 2016;113:E3101-10.

7    13. Cosgrove SE, Sakoulas G, Perencevich EN, Schwaber MJ, Karchmer AW, Carmeli Y.

8        Comparison of mortality associated with methicillin-resistant and methicillin-susceptible

9        Staphylococcus aureus bacteremia: a meta-analysis. Clin Infect Dis 2003;36:53–59.

10   14. Whitby M, McLaws ML, Berry G. Risk of death from methicillin-resistant Staphylococcus

11       aureus bacteraemia: a meta-analysis. Med J Aust. 2001;175:264 –267.

12   15. Melzer M, Eykyn SJ, Gransden WR, Chinn S. Is methicillin- resistant Staphylococcus

13       aureus more virulent than methicillin- susceptible S. aureus? A comparative cohort

14       study of British patients with nosocomial infection and bacteremia. Clin Infect Dis.

15       2003;37:1453–1460.

16   16. Laabei M, Recker M, Rudkin JK, Aldeljawi M, Gulay Z, Sloan TJ et al. Predicting the

17       virulence of MRSA from its genome sequence. Genome research, 2014;24:839-849.

18   17. Saunderson RB, Gouliouris T, Nickerson EK, Cartwright EJ, Kindey A, Aliyu SH, et al.

19       Impact of routine bedside infectious disease consultation on clinical management and

20       outcome of Staphylococcus aureus bacteraemia in adults. Clin Microbiol Infect.

21       2015;21:779-85.

22   18. Charlson ME, Pompei P, Ales KL, MacKenzie CR. A new method of classifying prognostic

23       comorbidity in longitudinal studies: Development and validation. J Chronic Dis.

24       1987;40:373–83.

1    19. Lesens O, Methlin C, Hansmann Y. Role of comorbidity in mortality related to

2        Staphylococcus aureus bacteremia: a prospective study using the Charlson weighted

3        index of comorbidity. Infect Control Hosp Epidemiol.2003;24:890–96.

4    20. Köser CU, Holden MT, Ellington MJ, Cartwright EJ, Brown NM, Ogilvy-Stuart AL, et al.

5        Rapid whole-genome sequencing for investigation of a neonatal MRSA outbreak. N Engl

6        J Med. 2012;366:2267-2275.

7    21. Holden MT, Hsu LY, Kurt K, Weinert LA, Mather AE, Harris SR, et al A genomic portrait of

8        the emergence, evolution, and global spread of a methicillin-resistant Staphylococcus

9        aureus pandemic. Genome Res. 2013;23:653-64.

10   22. Ziebuhr W, Krimmer V, Rachid S, Lossner I, Gotz F, Hacker J. A novel mechanism of phase

11       variation of virulence in Staphylococcus epidermidis: evidence for control of the

12       polysaccharide intercellular adhesin synthesis by alternating insertion and excision of

13       the insertion sequence element IS256. Mol Microbiol. 1999;32:345-56.

14   23. R Developement Core Team. R: A Language and Environment for Statistical Computing.

15       R Found Stat Comput. 2015;1:409.

16   24. Breiman L. Random forests. Mach Learn. 2001;45:5–32.

17   25. Liaw A, Wiener M. Classification and Regression by random forest. R News. 2002;2:18–

18       22.

19   26. Pozzi C, Waters EM, Rudkin JK, Schaeffer CR, Lohan AJ, et al. Methicillin resistance alters

20       the biofilm phenotype and attenuates virulence in Staphylococcus aureus device-

21       associated infections. PLoS Pathog. 2012;8(4):e1002626

22   27. Fey PD, Endres JL, Yajjala VK, Widhelm TJ, Boissy RJ, Bose JL, et al. A genetic resource for

23       rapid and comprehensive phenotype screening of nonessential Staphylococcus aureus

24       genes. MBio. 2013;4:e00537-12.
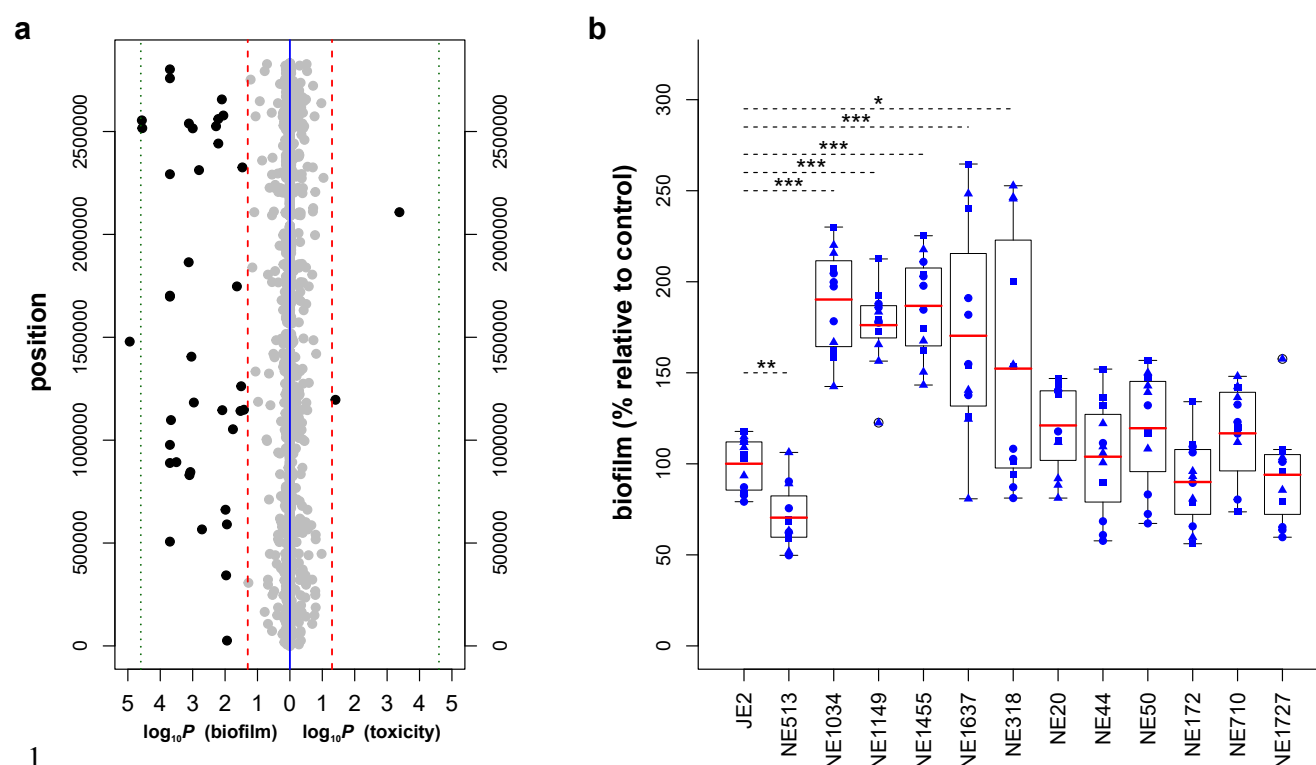
1  28. Nilsson IM, Lee JC, Bremell T, Rydén C, Tarkowski A. The role of staphylococcal

2      polysaccharide microcapsule expression in septicemia and septic arthritis. Infect Immun.

3      1997;65:4216-21.

4  29. Tzianabos AO, Wang JY, Lee JC. Structural rationale for the modulation of abscess

5      formation by Staphylococcus aureus capsular polysaccharides. Proc Natl Acad Sci USA.

6      2001;98:9365-70.

7  30. Boyle-Vavra S, Li X, Alam MT, Read TD, Sieth J, Cywes-Bentley C, et al. USA300 and

8      USA500 clonal lineages of Staphylococcus aureus do not produce a capsular

9      polysaccharide due to conserved mutations in the cap5 locus. MBio. 2015;6:e02585-14.

10

11

12

13

14

15

16

17

18

19

20

21

22

23

1 **Figures**



2
3

4 **Figure 1. The toxicity and biofilm forming abilities of 135 *S. aureus* bacteraemia isolates.**

5 **(a)** Toxicity for each isolate was determined by incubating bacterial supernatant with

6 cultured human cells, using flow cytotometry to quantify cell death (toxicity). No difference

7 was observed between methicillin-resistant (MRSA, grey circles) and methicillin-susceptible

8 (MSSA, blue circles) isolates. **(b)** Biofilm forming abilities were quantified relative to a control

9 included in each assay in a static 96 well format. A wide range of biofilm forming abilities

10 was evident with no discernible difference between methicillin-resistant (MRSA, grey circles)

11 and methicillin-susceptible (MSSA, blue circles) isolates.

12

13

**Figure 2. Genome-wide associations and functional validation of biofilm affecting polymorphisms. (a)** Manhattan plots representing the results of the GWAS for both biofilm (left-hand side) and toxicity (right-hand side), performed on the 135 bacteremia isolates. For toxicity, only two polymorphic loci were significantly associated using an uncorrected threshold (indicated by the red vertical dashed line), one in the *agrC* gene and the other in a gene encoding a putative membrane protein. No loci were significantly associated with toxicity when Bonferroni was used to correct for multiple comparisons (indicated by the green vertical dotted line). For biofilm, one locus, the *pbp2* gene, was significantly associated using the the Bonferroni threshold, and a f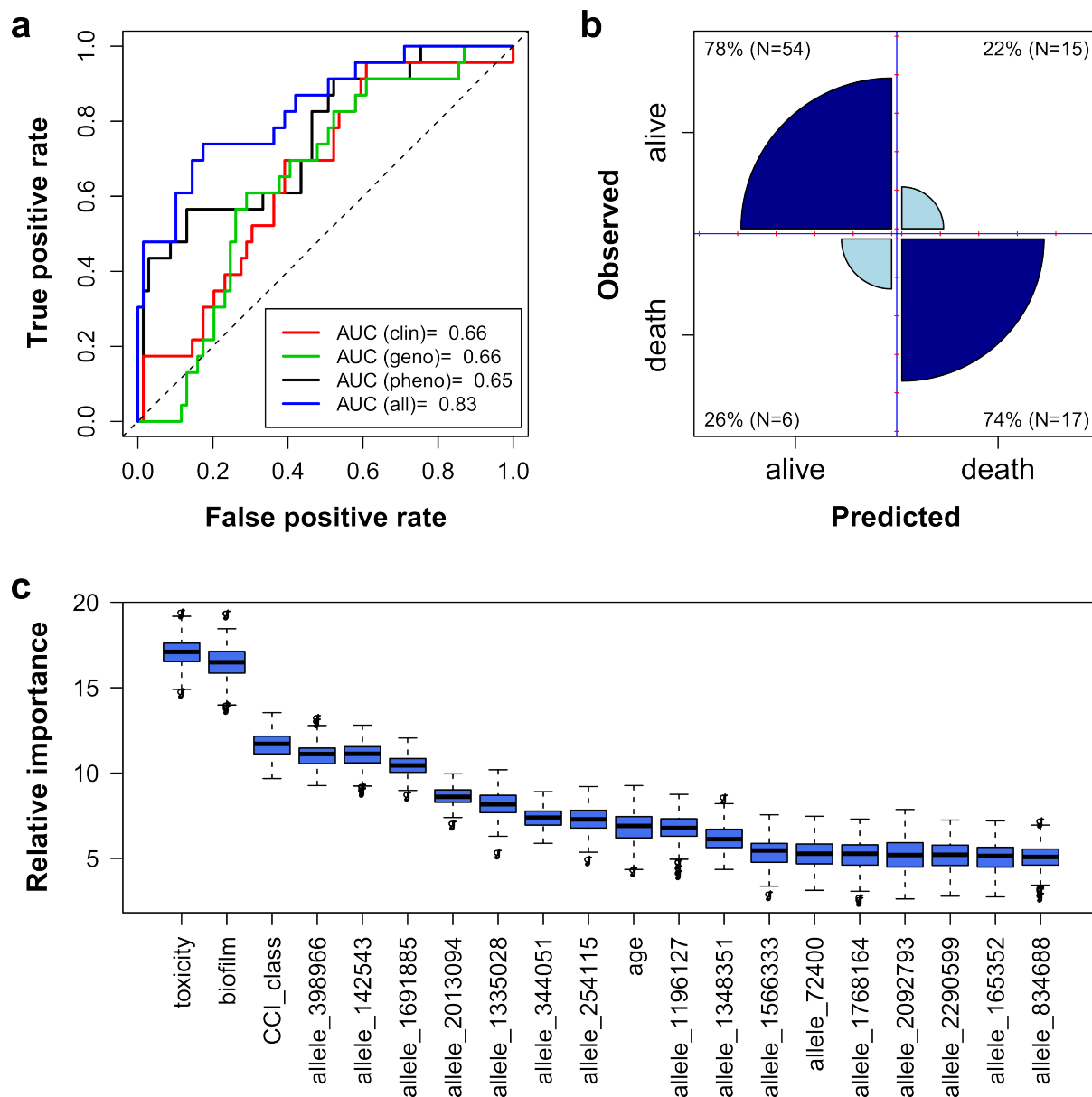urther 36 using the uncorrected threshold. **(b)** Twelve of the putative biofilm-affecting loci (associated by GWAS) were functionally validated using transposon insertion, of which six showed a significant effect compared to the wild-type (JE2) (Welch's two-sided *t*-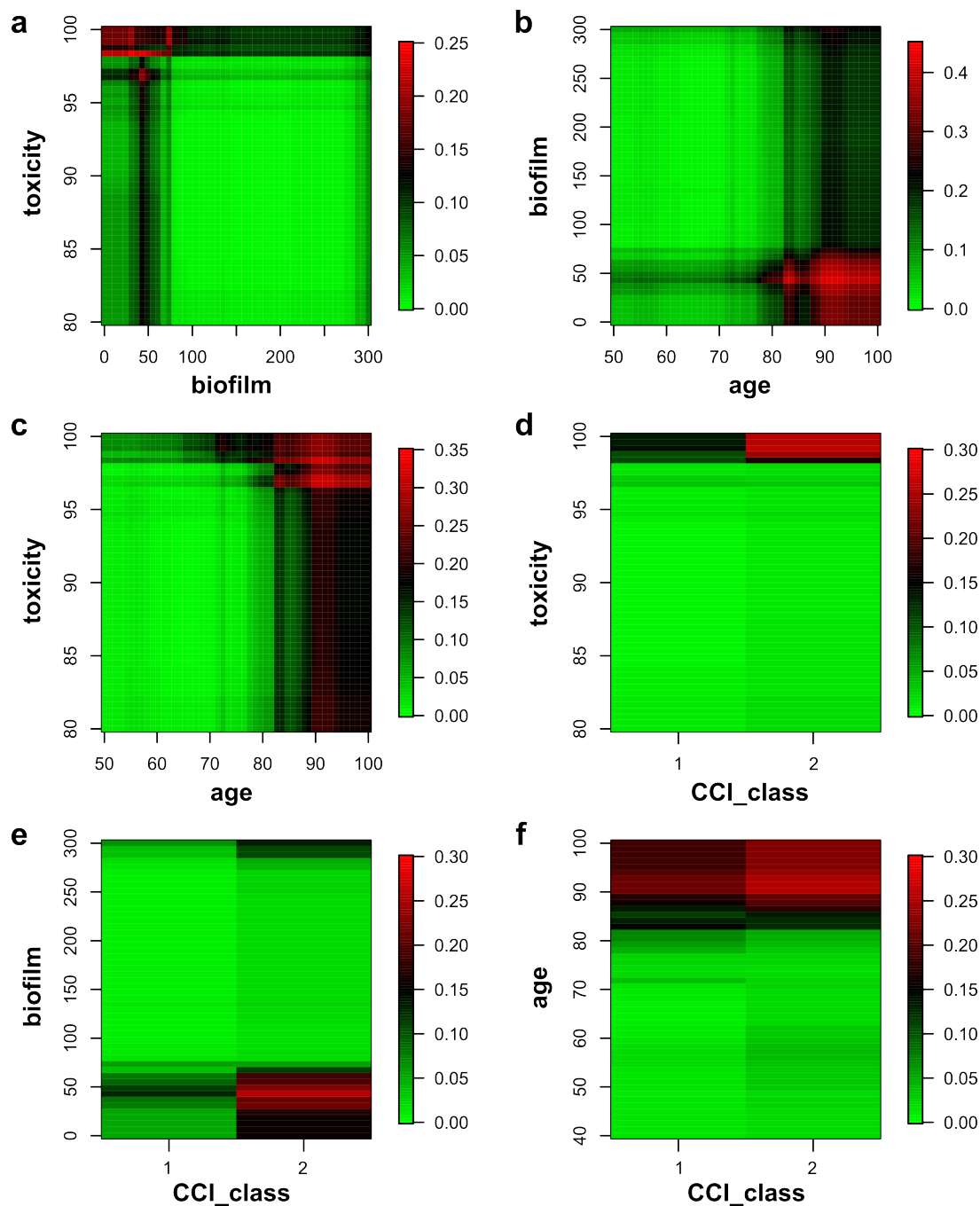test, with *: $P<2E-2$, **: $P<2E-4$, ***: $P<2E-6$). Results are based on four biological replicates per strain that were repeated three times each.

20

1

2 **Figure 3. Predictive model performance and variable importance. (a)** Receiver operator

3 characteristic (ROC) curve of the *random forests* fit to four sets of variables consisting of

4 clinical metadata (red line), genotype data (green line), phenotype data (black line) and all

5 available variables (blue lie). As indicated by the area under curve (AUC), the model fit to a

6 combination of all available data showed the highest predictive accuracy. **(b)** Confusion

7 matrix illustrating the accuracy of the model in predicting individual patients 30-day

8 mortality. The out-of-bag classification and misclassification rates (dark-blue/diagonal and

9 light-blue/off-diagonal wedges, respectively) of the *random forests* model fitted to a feature-

1    selected subset of the combined genotype, phenotype and clinical data. **(c)** Relative

2    importance of the 20 features used for the final model, as measured by the random forest

3    by means of a variable's influence on the model's predictive performance and based on 200

4    model fits, clearly demonstrating the significant effect of the bacteria's phenotype on SAB-

5    associated mortality.

6

Figure 4. Interactions of phenotype and host factors determine risk of mortality. Two-way interaction plots derived by the *random forests* visualising the interactions between the two phenotype measures (toxicity and biofilm) and two important host factors and their combined effect on host mortality. The color bars indicate the risk of mortality, keeping all other variables fixed at either their means (for continuous variable) or their most common value (for categorical predictors).

6   **Supporting Information Captions.**

7   **Supplementary Table 1:** Strain information including clinical metadata, genome accession
8   codes, toxicity and biofilm formation. (this has been provided as an excel spreadsheet)
9
10  **Supplementary Table 2:** Loci associated by GWAS with biofilm formation.

| SNP position | Locus tag in the *S. aureus* reference strain HO50960412 | Gene name (if ascribed) | Putative/known Activity |
|---|---|---|---|
| 25928 | SAEMRSA15_0190 | *yycG* | Sensor kinase protein |
| 342920 | SAEMRSA15_02780 | | NADH oxidase family protein |
| 506893 | Intergenic between SAEMRSA15_04340 and SAEMRSA15_04350 | | |
| 566273 | SAEMRSA15_04760 | | Putative peptidase |
| 590633 | SAEMRSA15_04910 | | Putative gyclosyltransferase |
| 661898 | SAEMRSA15_05660 | | Putative gyclosyltransferase |
| 830176 | SAEMRSA15_07180 | *nuc* | Thermonuclease precursor |
| 844980 | SAEMRSA15_07420 | | Putative lipoprotein |
| 889042 | SAEMRSA15_07880 | | Putative NAD-specific glutamate dehydrogenase |
| 893074 | SAEMRSA15_07910 | *argG* | Putative argininosuccinate synthethase |
| 977070 | SAEMRSA15_08640 | | pseudogene |
| 1053036 | Intergenic between SAEMRSA15_09350 and SAEMRSA15_09360 | | |
| 1097556 | SAEMRSA15_09780 | *sdhA* | Putative succinate dehydrogenase flavoprotein subunit |
| 1142055 | SAEMRSA15_10260 | *ileS* | Isoleucyl tRNA synthetase |
| 1145714 | SAEMRSA15_10300 | | Putative RNA pseudouridylate |

| | | | synthase |
|---|---|---|---|
| 1147728 | SAEMRSA15_10320 | *pyrP* | Putative uracil permease |
| 1182957 | SAEMRSA15_10630 | *fabD* | Putative malonyl CoA-acyl carrier protein |
| 1262036 | SAEMRSA15_11300 | mutS | DNA mismatch repair |
| 1405808 | SAEMRSA15_12710 | | Conserved hypothetical protein |
| 1479055 | SAEMRSA15_13110 | *pbp2* | Penicillin binding protein |
| 1697502 | SAEMRSA15_15450 | | Putative ATPase |
| 1702225 | SAEMRSA15_15480 | *hisS* | Histidyl tRNA synthetase |
| 1747240 | SAEMRSA15_15930 | *thrS* | Threonyl tRNA synthetase |
| 1864667 | SAEMRSA15-16820 | *norA* | quinolone efflux protein |
| 2292604 | Intergenic between SAEMRSA15_21020 and SAEMRSA15_21030 | | |
| 2312449 | SAEMRSA15_21210 | | ABC transporter ATP-binding protein |
| 2325657 | SAEMRSA15_21480 | *rplD* | 50S ribosomal protein L4 |
| 2441578 | SAEMRSA15_22660 | *mqo1* | Putative malate:quinone oxidoreductase |
| 2515404 | SAEMRSA15_23310 | | Putative membrane protein |
| 2516262 | SAEMRSA15_23320 | | Putative glycerate kinase |
| 2525636 | SAEMRSA15_23400 | | Putative membrane protein |
| 2538962 | SAEMRSA15_23500 | | Putative amino acid permease |
| 2554304 | SAEMRSA15_23640 | *opp-1B* | Putative oligopeptide transporter membrane protein |
| 2560821 | SAEMRSA15_23690 | | Hypothetical protein |
| 2577780 | SAEMRSA15_23880 | | Putative helicase |
| 2656043 | SAEMRSA15_24600 | *copA* | Putative copper importing ATPase |
| 2758770 | SAEMRSA15_25490 | | Putative exported protein |
| 2802186 | SAEMRSA15_25850 | *hisG* | Putative ATP phosphoribosyltransferase |

1
2
3
4 **Supplementary Table 3:** Bacterial loci identified by the model as important in predicting SAB-
5 associated mortality (fig. 4c).
6
7

| Allele | Locus tag in the *S. aureus* reference strain HO50960412 | Gene name (if ascribed) | Putative/known function |
|---|---|---|---|
| 398966 | SAEMRSA15_03350 | | nitroreductase family protein |
| 142543 | SAEMRSA15_01150 | *capA* | capsular polysaccharide synthesis enzyme |

| 1691885 | intergenic between SAEMRSA15_15380 and 15390 | | |
|---|---|---|---|
| 2013094 | SAEMRSA15_18250 | | putative prephenate dehydratase |
| 1335028 | SAEMRSA15_12030 | *grlA* | topoisomerase IV subunit A |
| 344051 | SAEMRSA15_02790 | | putative bacterial luciferase family protein |
| 254115 | SAEMRSA15_02010 | | putative PTS transport system, IIBC component |
| 1196127 | SAEMRSA15_10750 | | putative membrane protein |
| 1348351 | SAEMRSA15_12130 | | prephenate dehydrogenase |
| 1566333 | intergenic between SAEMRSA15_14090 and 14100 | | |
| 72400 | SAEMRSA15_00580 | | LysR-family regulatory protein |
| 1768164 | SAEMRSA15_16080 | *pfkA* | 6-phosphofructokinase |
| 2092793 | intergenic between SAEMRSA15_19290 and 19300 | | |
| 165352 | SAEMRSA15_01380 | | putative lipoprotein |
| 834688 | SAEMRSA15_07270 | | phosphoglycerate mutase family protein |

1

2

3

4

5

6

7

8

9

10