

1 **Panton-Valentine leucocidin is the key determinant of *Staphylococcus aureus***
2 **pyomyositis in a bacterial genome-wide association study**

3 Bernadette C Young^{1,2}, Sarah G Earle¹, Sona Soeng³, Poda Sar³, Varun Kumar⁴, Songly Hor³,
4 Vuthy Sar³, Rachel Bousfield⁵, Nicholas D Sanderson¹, Leanne Barker¹, Nicole Stoesser^{1,6},
5 Katherine RW Emary², Christopher M Parry^{7,8}, Emma K Nickerson⁵, Paul Turner^{3,9}, Rory
6 Bowden¹⁰, Derrick Crook^{1,2,6}, David Wyllie^{1,6,11}, Nicholas PJ Day^{9,13}†, Daniel J Wilson^{1,10,12} †,
7 Catrin E Moore^{9,13}†*

8 ¹ Nuffield Department of Medicine, Experimental Medicine Division, University of Oxford,
9 John Radcliffe Hospital, Oxford OX3 9DU, UK.

10 ² NIHR Oxford Biomedical Research Centre, Infection Theme, Oxford University Hospitals
11 NHS Foundation Trust, John Radcliffe Hospital, Oxford, OX3 9DU, UK.

12 ³ Cambodia Oxford Medical Research Unit, Angkor Hospital for Children, Siem Reap,
13 Cambodia

14 ⁴ Department of Pediatrics, East Tennessee State University Quillen College of Medicine,
15 Johnson City, USA

16 ⁵ Department of Infectious Diseases, Cambridge University Hospitals NHS Foundation Trust,
17 Cambridge, CB2 0QQ, UK

18 ⁶ Public Health England Academic Collaborating Centre, John Radcliffe Hospital, Oxford OX3
19 9DU, UK.

20 ⁷ Clinical Sciences, Liverpool School of Tropical Medicine, Liverpool, L3 5QA

21 ⁸ School of Tropical Medicine and Global Health, Nagasaki University, Nagasaki, Japan

22 ⁹ Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine,
23 University of Oxford, Oxford OX3 7ZF, UK

24 ¹⁰ Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK

25 ¹¹ The Jenner Institute Laboratories, University of Oxford, Old Road Campus Research
26 Building, Roosevelt Drive, Oxford OX3 7DQ, UK.

27 ¹² Institute for Emerging Infections, Oxford Martin School, University of Oxford, Oxford, OX1
28 3BD, UK

29 ¹³ Mahidol-Oxford Tropical Medicine Research Unit, Faculty of Tropical Medicine, Mahidol
30 University, Bangkok 10400, Thailand

31

32 *Correspondence to: catrin.moore@ndm.ox.ac.uk

33 †These authors contributed equally to the work

34

35 **Abstract:** Pyomyositis is a severe bacterial infection of skeletal muscle, commonly affecting
36 children in tropical regions and predominantly caused by *Staphylococcus aureus*. To
37 understand the contribution of bacterial genomic factors to pyomyositis, we conducted a
38 genome-wide association study of *S. aureus* cultured from 101 children with pyomyositis and
39 417 children with asymptomatic nasal carriage attending the Angkor Hospital for Children in
40 Cambodia. We found a strong relationship between bacterial genetic variation and pyomyositis,
41 with estimated heritability 63.8% (95% CI 49.2-78.4%). The presence of the Pantan-Valentine
42 leucocidin (PVL) locus increased the odds of pyomyositis 130-fold ($p=10^{-17.9}$). The signal of
43 association mapped both to the PVL-coding sequence and the sequence immediately upstream.
44 Together these regions explained >99.9% of heritability. Our results establish staphylococcal
45 pyomyositis, like tetanus and diphtheria, as critically dependent on expression of a single toxin
46 and demonstrate the potential for association studies to identify specific bacterial genes
47 promoting severe human disease.
48
49

50 Introduction

51 Microbial genome sequencing and bacterial genome-wide association studies present new
52 opportunities to discover bacterial genes involved in the pathogenesis of serious infections.¹⁻⁶
53 Pyomyositis is a severe infection of skeletal muscle most commonly seen in children in the
54 tropics.⁷⁻⁹ In up to 90% of cases it is caused by a single bacterial pathogen, *Staphylococcus*
55 *aureus*.⁷⁻¹⁰ There is evidence that some *S. aureus* strains have heightened propensity to cause
56 pyomyositis – the incidence in the USA doubled during an epidemic of community-associated
57 methicillin resistant *S. aureus* (CA-MRSA)¹¹ – but molecular genetic investigation of *S. aureus*
58 from pyomyositis has been limited.¹²

59 Panton-Valentine leucocidin (PVL), a well-known staphylococcal toxin causing purulent skin
60 infections and found in epidemics caused by CA-MRSA, has been implicated in pyomyositis,
61 pneumonia and other *S. aureus* disease manifestations, but its role is strongly disputed.¹³⁻¹⁶
62 PVL is a bipartite pore-forming toxin comprising the co-expressed LukF-PV and LukS-PV
63 proteins,^{17,18} is encoded by *lukSF-PV*, which is usually carried on bacteriophages^{13,17} which
64 facilitates *lukSF-PV* exchange between lineages.¹⁹ Although small case series testing for
65 candidate genes have reported a high prevalence of PVL among pyomyositis-causing
66 *S. aureus*,^{11,20,21} a detailed meta-analysis found no evidence for an increased rate of
67 musculoskeletal infection or other invasive disease in PVL-positive bacteria *versus* controls.¹³
68 These conflicting results may reflect insufficiently powered studies. However, candidate gene
69 studies may also miss important variation elsewhere in the genome: a study claiming a critical
70 role for PVL in the causation of severe pneumonia¹⁵ was later found to have overlooked
71 mutations in key regulatory genes, capable of producing the virulent behaviour that had been
72 attributed to PVL by the original study.¹⁶ Thus, based on the current evidence, opinion is
73 divided as to whether PVL is an important virulence factor in pyomyositis, or merely an
74 epiphenomenon, carried by bacteria alongside unidentified genetic determinants.^{22,23}

75 Genome-wide association studies (GWAS) offer a means to screen entire bacterial genomes to
76 discover genes and genetic variants associated with disease risk. They are particularly
77 appealing because they enable the investigation of traits not readily studied in the laboratory,
78 and do not require the nomination of specific candidate genes.⁵ Proof-of-principle GWAS in
79 bacteria have already demonstrated the ability to successfully rediscover known antimicrobial
80 resistance (AMR) determinants.^{2,3,4} However, AMR is under extraordinarily intense selection in
81 bacteria, and it remains to be seen whether GWAS can overcome the typically strong linkage
82 disequilibrium in bacterial populations to precisely pinpoint genes and genetic variants
83 underlying the propensity to cause human infection.

84 Results

85 To understand the bacterial genetic basis of pyomyositis, we sampled and whole-genome
86 sequenced *S. aureus* from 101 pyomyositis infections and 417 asymptomatic nasal carriage
87 episodes in 518 children attending Angkor Hospital for Children in Siem Reap, Cambodia
88 between 2008 and 2012 (Table S1). As expected of *S. aureus* epidemiology, we observed
89 representatives of multiple globally common lineages in Cambodia, together with some
90 globally less common lineages at high frequency, in particular clonal complex (CC) 121,
91 identified by Multi-locus sequence type (MLST). Lineage composition appeared stable over
92 time, with no major changes in lineage frequency (Fig. S1).

93 In our study, some *S. aureus* lineages were strongly overrepresented among cases of
94 pyomyositis compared to their frequency among asymptomatic, nasally-carried controls over
95 the same time period. Notably, 86/101 (85%) of pyomyositis cases were caused by CC-121
96 bacteria, whereas no pyomyositis cases were caused by the next two most commonly carried

97 lineages, sequence type (ST)-834 and CC-45 (Fig. 1). We estimated the overall heritability of
98 case/control status to be 63.8% (95% CI 49.2-78.4%) in the sample, reflecting the strong
99 relationship between bacterial genetic variation and case/control status. We used *bugwas*⁶ to
100 decompose this heritability into the principal components (PCs) of bacterial genetic variation.
101 PC 1, which distinguished CC-121 (the most common pyomyositis lineage) from ST-834
102 (which was only found in carriage), showed the strongest association with case/control status
103 ($p = 10^{-29.6}$, Wald test). The next strongest were PC 20, which differentiated a sub-lineage of
104 CC-121 within which no cases were seen ($p = 10^{-13.9}$), and PC 2, which distinguished CC-45
105 from the rest of the species ($p = 10^{-4.9}$).

106 We conducted a GWAS to identify bacterial genetic variants associated with pyomyositis,
107 controlling for differences in pyomyositis prevalence between *S. aureus* lineages. We used a
108 kmer-based approach¹ in which every variably present 31bp DNA sequence observed among
109 the 518 genomes was tested for association with pyomyositis *versus* asymptomatic nasal
110 carriage, controlling for population structure using GEMMA.²⁴ These kmers captured bacterial
111 genetic variation including single nucleotide polymorphisms (SNPs), insertions or deletions
112 (indels), and presence or absence of entire accessory genes. We found 10.7 million unique
113 kmers variably present across the bacterial genomes. In total, 9,175 kmers were significantly
114 associated with case/control status after correction for multiple testing ($10^{-6.8} \leq p \leq 10^{-21.4}$; Fig
115 2A). The vast majority of these kmers (8,993/9,175; 98.0%) localised to a 45.7kb region
116 spanning an integrated prophage with 95% nucleotide similarity to ϕ SLT (Fig 2B). Most
117 significant kmers, (9,173/9,175; 99.98%) were found more frequently in pyomyositis, with
118 odds ratios (OR) ranging from 2.7 to 139.8, indicating they were associated with increased risk
119 of disease. The bacteriophage ϕ SLT was thus strongly associated with pyomyositis.

120 We were able to fine-map the signal of association within ϕ SLT to the *lukS-PV* and *lukF-PV*
121 cargo genes. These genes encode the subunits of PVL, which multimerise into a pore-forming
122 toxin capable of rapidly lysing the membranes of human neutrophils, the first line of defense
123 against *S. aureus*.^{17,25} 1630 kmers tagging the presence of the *lukSF-PV* coding sequences
124 (CDS) were highly significantly associated with disease, being present in 98/101 (97%)
125 pyomyositis cases and 84/417 (20%) carriage controls (OR 129.5, $p=10^{-17.9}$). Kmers tagging
126 variation in the 389bp region immediately upstream of the CDS were also strongly associated
127 with disease ($p=10^{-21.4}$). The most significant of these kmers were co-present with the CDS in
128 the same cases (98/101, 97.0%), but present in fewer controls (79/417, 18.9%), producing an
129 OR of 140.

130 Closer examination of this ~400bp upstream region in genomes assembled from short-read
131 Illumina sequencing showed that assembly of the region was problematic, with breaks or gaps
132 in the assembly (Fig. S2). To improve the accuracy of this region of the assembled genomes we
133 performed long-read Oxford Nanopore sequencing on the 37 genomes with incomplete or
134 discontinuous assembly upstream of the PVL CDS. By integrating long-read and short-read
135 data we were able to assemble a single contig spanning this region in all isolates (Fig. S3).
136 When these improved assemblies were introduced, the signal of association upstream of the
137 PVL CDS was no more significant than within the CDS (Fig 2C). Therefore, the presence of
138 genomic sequence spanning the PVL toxin-coding sequences and the upstream, presumed
139 regulatory, region exhibited the strongest association with pyomyositis in the *S. aureus*
140 genome.

141 Presence or absence of the PVL region accounted for the differences in pyomyositis rates
142 between lineages. It was common in pyomyositis-associated lineages including CC-121 and
143 absent from non-pyomyositis-associated lineages including ST-834 and CC-45 (Fig. 1),
144 explaining 99.9% of observed heritability in case-control status. It was infrequent in the non-

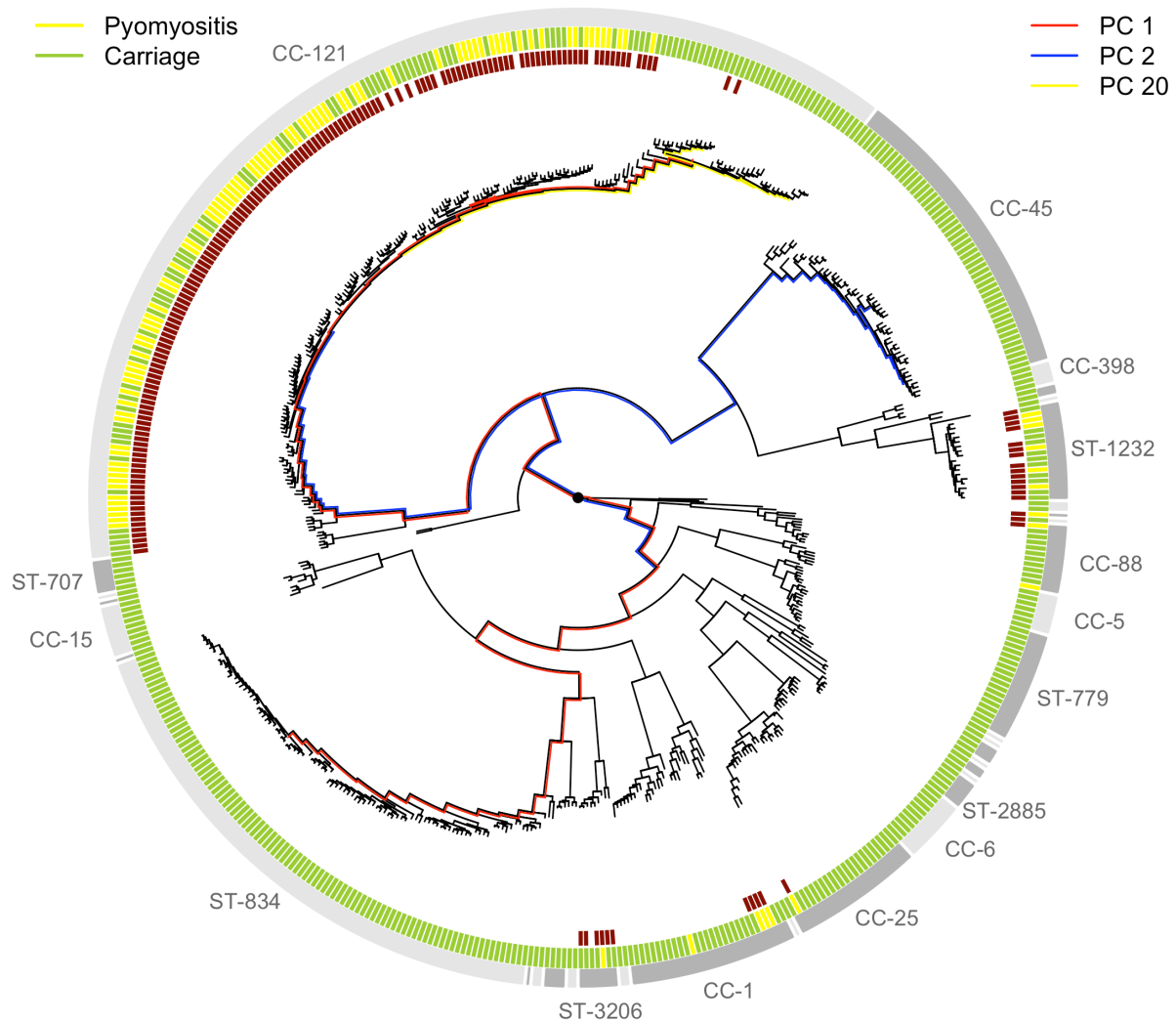
145 pyomyositis-associated sub-lineage of CC-121 (2/36, 5.6%), and sporadically present in
146 pyomyositis cases in otherwise non-pyomyositis-associated, PVL-negative strains CC-1 and
147 CC-88. Its absence from only three cases (in CC-88, CC-1 and CC-121) suggested that the PVL
148 region approached necessity for development of pyomyositis in Cambodian children, while its
149 presence in 20% of controls indicated that PVL-associated pyomyositis is incompletely
150 penetrant, i.e. presence of the PVL region does not always lead to disease.

151 **Discussion**

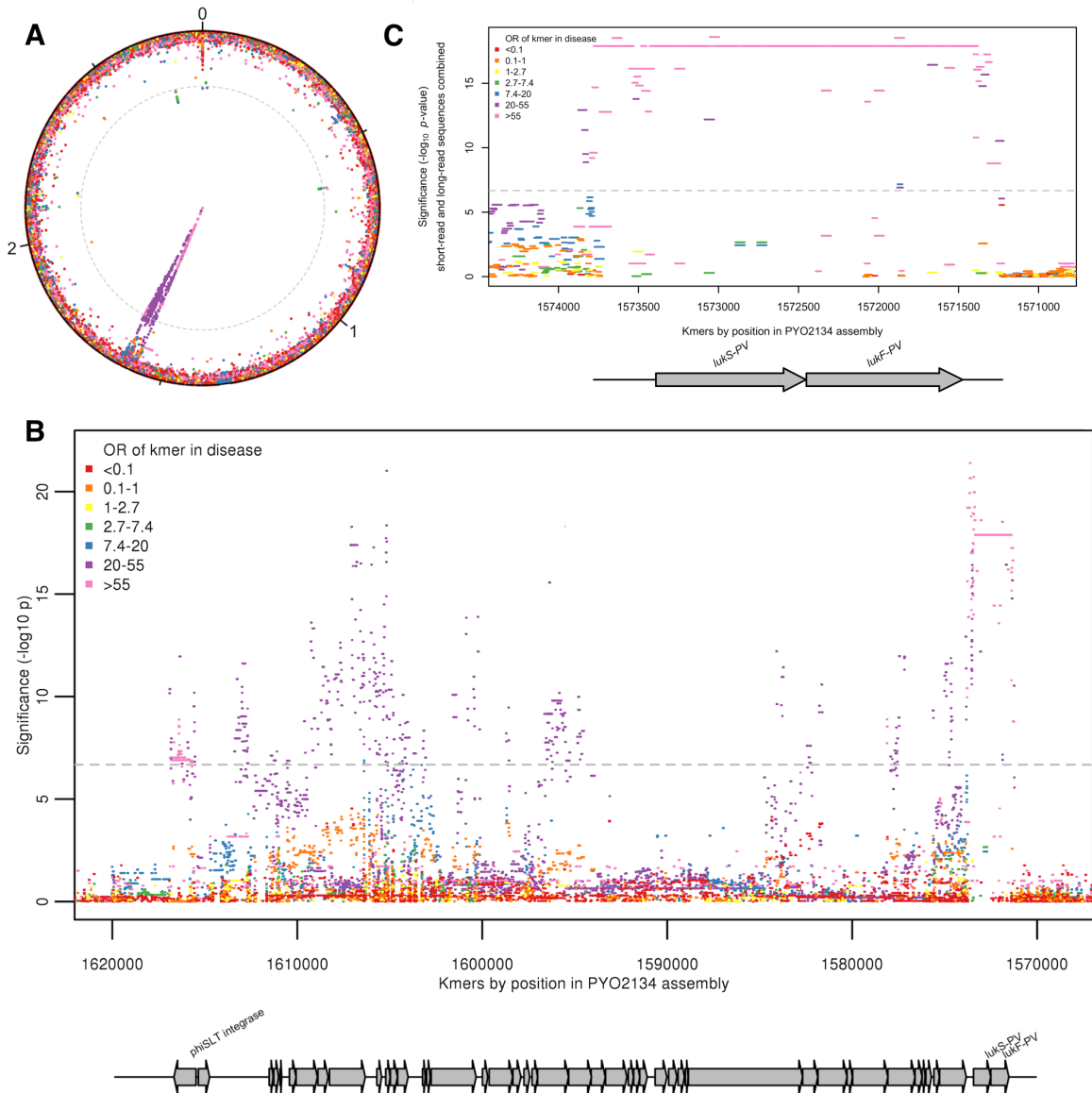
152 In this study we found a strong association between pyomyositis, a highly distinctive tropical
153 infection of skeletal muscle in children, and Pantone-Valentine leukocidin, a bacterial toxin
154 commonly carried by bacteriophages. We found that a single coding region together with the
155 upstream sequence are all but necessary for the development of pyomyositis: its sporadic
156 presence is associated with pyomyositis in otherwise low-frequency strains, and its absence is
157 associated with asymptomatic carriage in a high propensity strain. The locally common PVL-
158 positive CC-121 lineage contributes most strongly to the prevalence of pyomyositis in
159 Cambodian children.

160 While PVL has long been thought an important *S. aureus* virulence factor,^{25,26,27} its role in
161 invasive disease has been controversial,^{22,23} with conflicting results in case-control studies and
162 an absence of supporting evidence on meta-analysis.¹³ In previous studies the PVL positive
163 USA300 lineage was associated with musculoskeletal infection (both pyomyositis and
164 osteomyelitis), however in these studies almost all these infections were caused by the USA300
165 strain, so the role of PVL was almost completely confounded by both methicillin resistance and
166 strain background.^{11,26} In our study, this confounding is broken down by the movement of PVL
167 on mobile genetic elements (MGEs). Here, despite the emergence of CA-MRSA in carriage in
168 the same population,²⁸ all the pyomyositis cases were MSSA. By applying GWAS methods to a
169 well-powered cohort, our study resolves the controversy around PVL and pyomyositis,
170 demonstrating strong heritability which localises to a single region, even when the full bacterial
171 genome is considered. This study offers empirical evidence that, in addition to elucidating
172 phenotypes under strong selection (such as in antimicrobial resistance)^{2,3,4}, bacterial GWAS can
173 pinpoint variants when MGEs act to unravel linkage disequilibrium.

174 These results offer a novel prospect for disease prevention: they establish staphylococcal
175 pyomyositis as a disease whose pathogenesis depends critically on expression of a single toxin.
176 This property is shared by toxin-driven, vaccine-preventable diseases such as tetanus and
177 diphtheria. Therefore, vaccines that generate neutralising anti-toxin antibodies against PVL²⁹
178 may protect human populations against this common tropical disease. They also raise the
179 hypothesis that antibiotics which decrease toxin expression, and have been recommended in
180 some PVL-associated infections,³⁰ may offer specific clinical benefit in treating pyomyositis.
181 More generally, our study provides an example of how microbial GWAS can be used to
182 elucidate the pathogenesis of bacterial infections and identify specific virulence genes
183 associated with human disease.



184
185 **Figure 1.** Phylogeny of *S. aureus* cultured from children in Cambodia shows strong strain-to-
186 strain variation in pyomyositis prevalence. The phylogeny was estimated by maximum
187 likelihood from SNPs mapping to the USA300 FPR3757 reference genome. Multi-locus
188 sequence type (ST) or clonal complex (CC) groups are shown (outer gray ring). Case/control
189 status is marked in the middle ring: pyomyositis (gold, n = 101) or nasal carriage (green, n =
190 417). Branches of the phylogeny that correspond to the three principal components (PCs)
191 significantly associated with case/control status (PCs 1, 2 and 20) are marked in red, blue and
192 yellow, respectively. Branch lengths are square root transformed to aid visualization. The
193 presence of the kmers most strongly associated with pyomyositis is indicated by red blocks in
194 the inner ring



195

196 **Figure 2.** Kmers associated with pyomyositis. (A) All kmers (n = 10,744,013) were mapped to
 197 the genome assembly of one CC121 pyomyositis bacterium (PYO2134). Each point represents
 198 a kmer, plotted by the mapped location and the significance of the association with disease (-
 199 log₁₀ p value). Kmers are coloured by the odds ratio (OR) of kmer presence for disease risk. A
 200 Bonferroni-adjusted threshold for significance is plotted in grey (B) The region between 1.57-
 201 1.62 MB in greater detail. Gray arrows depict coding sequences, determined by homology to
 202 USA300 FPR3757. (C) Associations for kmers mapping to region 1,571 – 1,574kB is plotted.
 203 Kmer presence determined from hybrid assembly using short and long-read data for assembly.
 204 Gray arrows depict coding sequences, determined by homology to USA300 FPR3757.

205

206 **Materials and Methods:**

207 **Patient sample collection.** We retrospectively identified pyomyositis cases from the Angkor
208 Hospital for Children in Siem Reap, Cambodia, between January 2007 and November 2011.
209 We screened all attendances in children (≤ 16 years) using clinical coding (ICD-10 code M60
210 (myositis)) and isolation of *S. aureus* from skeletal muscle abscess pus. We reviewed clinical
211 notes to confirm a clinical diagnosis of pyomyositis was made by the medical staff, and
212 bacterial strains cultured by routine clinical microbiology laboratory were retrieved from the
213 local microbiology biobank. 106 clinical episodes of pyomyositis were identified, in 101
214 individuals, and we included the earliest episode from each individual.

215 We identified *S. aureus* nasal colonisation from two cohort studies undertaken at Angkor
216 Hospital for Children. The first were selected from a collection characterising nasal
217 colonization in the region between September and October, 2008, which has previously been
218 described using multi-locus sequence typing.²⁸ The swabs had been saved at -80°C since the
219 study, these samples were reexamined for the presence of *S. aureus* using selective agar,
220 confirmed using Staphaurex (Remel, Lenexa, USA) and the DNase agar test (Oxoid,
221 Hampshire, UK). Antimicrobial susceptibility testing was performed according to the 2014
222 Clinical and Laboratory Standards Institute guidelines (M100-24).³¹

223 We undertook a second cohort study in 2012. Inclusion criteria were children (≤ 16 years)
224 attending as an outpatient at Angkor Hospital for Children with informed consent. There were
225 no exclusion criteria. Children were swabbed between the 2-7th July 2012, using sterile cotton
226 tipped swabs pre-moistened (with phosphate buffered saline, PBS) using 3 full rotations of the
227 swab within the anterior portion of each nostril with one swab being used for both nostrils, the
228 ends were broken into bottles containing sterile PBS and kept cool until plated in the laboratory
229 (within the hour). The swabs were plated onto Mannitol Salt agar to select for *S. aureus*. The
230 M100-24 CLSI³¹ standards were followed for susceptibility testing and bacteria stored in
231 tryptone soya broth and glycerol at -80°C .

232 We selected controls from carriers in these two cohorts using the excel randomization function:
233 222 of 519 from the 2008 cohort and 195 of 261 from the 2012 cohort.

234 **Ethical Framework.** Approval for this study was provided by the AHC institutional review
235 board and the Oxford Tropical Ethics Committee (507-12).

236 **Whole genome sequencing.** For each bacterial culture, a single colony was sub-cultured and
237 DNA was extracted from the sub-cultured plate using a mechanical lysis step (FastPrep;
238 MPBiomedicals, Santa Ana, CA) followed by a commercial kit (QuickGene; Autogen Inc,
239 Holliston, MA), and sequenced at the Wellcome Trust Centre for Human Genetics, Oxford on
240 the Illumina (San Diego, California, USA) HiSeq 2500 platform, with paired-end reads 150
241 base pairs long.

242 A subset of samples were sequenced using long-read sequencing technology. We selected 37
243 isolates with incomplete assembly upstream of the PVL locus, 22 with ambiguous base calls in
244 the assembly, and 15 where the region was assembled over 2 contigs. DNA was extracted using
245 Genomic Tip 100/G (Qiagen, Manchester, UK) and DNA libraries prepared using Oxford
246 Nanopore Technologies (ONT) SQK-LSK108 library kit (ONT, Oxford, UK) according to
247 manufacturer instructions. These were then sequenced on ONT GridION device integrated with
248 a FLO-MIN106 flow cell (ONT, Oxford, UK). ONT base calling was performed using Guppy
249 v.1.6.

250 **Variant calling.** For short-read sequencing we used Velvet³² v1.0.18 to assemble reads into
251 contigs *de novo*. Velvet Optimiser v2.1.7 was used to choose the kmer lengths on a per
252 sequence basis. The median kmer length was 123bp (IQR 119-123). To obtain multilocus

253 sequence types we used BLAST³³ to find the relevant loci, and looked up the nucleotide
254 sequences in the online database at <http://saureus.mlst.net/>. Strains that shared 6 of 7 MLST
255 loci were considered to be in the same Clonal Complex. Antibiotic sensitivity was predicted by
256 interrogating the assemblies for a panel of resistance determinants as previously described.³⁴

257 We used Stampy³⁵ v1.0.22 to map reads against reference genomes (USA300_FPR3757,
258 Genbank accession number CP000255.1).³⁶ Repetitive regions, defined by BLAST³³
259 comparison of the reference genome against itself, were masked prior to variant calling. Bases
260 were called at each position using previously described quality filters.³⁷⁻³⁹

261 After filtering ONT reads with filtlong v.0.2.0 (with settings filtlong -- min_length 1000 --
262 keep_percent 90 --target_bases 500000000 --trim --split 500), hybrid assembly of short
263 (Illumina) and long (ONT) reads were made, using Unicycler v0.4.5⁴⁰ (default settings). The
264 workflow for these assemblies is available at
265 <https://gitlab.com/ModernisingMedicalMicrobiology/MOHAWK>)

266 **Reconstructing the phylogenetic tree.** We constructed a maximum likelihood phylogeny of
267 mapped genomes for visualization using RAxML⁴¹ assuming a general time reversible (GTR)
268 model. To overcome a limitation in the presence of divergent sequences whereby RAxML fixes
269 a minimum branch length that may be longer than a single substitution event, we fine-tuned the
270 estimates of branch lengths using ClonalFrameML.⁴²

271 **Kmer counting.** We used a kmer-based approach to capture non-SNP variation.¹ Using the *de*
272 *novo* assembled genome, all unique 31 base haplotypes were counted using dsk⁴³. If a kmer
273 was found in the assembly it was counted present for that genome, otherwise it was treated as
274 absent. This produced a set of 10,744,013 variably present kmers, with the presence or absence
275 of each determined per isolate. We identified a median of 2,801,000 kmers per isolate,
276 including variably present kmers and kmers common to all genomes (IQR 2,778,000-
277 2,837,000).

278 **Calculating heritability.** We used the Genome-wide Efficient Mixed Model Association tool
279 (GEMMA²⁴) to fit a univariate linear mixed model for association between a single phenotype
280 (pyomyositis vs asymptomatic nasal carriage). We calculated the relatedness matrix from
281 kmers, and used GEMMA to estimate the proportion of variance in phenotypes explained by
282 genotypic diversity (i.e. heritability).

283 **Genome wide association testing of Kmers.** We performed association testing using an R
284 package bacterialGWAS (<https://github.com/jessiewu/bacterialGWAS>), which implements a
285 published method for locus testing in bacterial GWAS.³ The association of each kmer on the
286 phenotype was tested, controlling for the population structure and genetic background using a
287 linear mixed model (LMM) implemented in GEMMA.²⁴ The parameters of the linear mixed
288 model were estimated by maximum likelihood and likelihood ratio test was performed against
289 the null hypothesis (that each locus has no effect) using the software GEMMA.²⁴ GEMMA was
290 run using a minor allele frequency of 0 to include all SNPs. GEMMA was modified to output
291 the ML log-likelihood under the null, and alternative and $-\log_{10} p$ values were calculated using
292 R scripts in the bacterialGWAS package. Unadjusted odds ratios were reported because there
293 was no residual heritability unexplained by the most significant kmers.

294 **Testing for lineage effects.** We tested for associations between lineage and phenotype using an
295 R package *bugwas* (available at <https://github.com/sgearle/bugwas>), which implements a
296 published method for lineage testing in bacterial GWAS.³ We tested lineages using principal
297 components. These were computed based on biallelic SNPs using the R function *prcomp*. To
298 test the null hypothesis of no background effect of each principal component, we used a Wald

299 test, which we compared against a χ^2 distribution with one degree of freedom to obtain a p
300 value.

301 **Kmer mapping.** We used Bowtie⁴⁴ to align all 31bp kmers from short-read sequencing were to
302 a draft reference (the *de novo* assembly of a CC-121 pyomyositis strain PYO2134). Areas of
303 homology between the draft reference and well-annotated reference strains were identified by
304 aligning sequences with Mauve⁴⁵. For all 31bp kmers with significant association with case-
305 controls status, the likely origin of the kmer was determined by nucleotide sequence BLAST³³
306 of the kmers against a database of all *S. aureus* sequences in Genbank.

307 **Joint short-read and long-read analysis.** 31bp kmers were counted for the 37 hybrid short-read
308 and long-read assemblies using dsk⁴³. The presence or absence of all Illumina (short-read)
309 kmers that mapped to the two PVL toxin-coding sequences and the upstream intergenic region
310 plus the surrounding 1kb were reassessed. For the 37 samples with hybrid assemblies, the
311 presence/absence of these kmers was determined from the kmers counted from the hybrid
312 assemblies. For all other samples, presence/absence was determined from the kmers counted
313 from the short-read only assemblies. The new presence/absence patterns were tested for
314 association with the phenotype controlling for population structure and genetic background
315 using GEMMA²⁴, using the same relatedness matrix as the original short-read analysis.

316 **Multiple testing correction.** Multiple testing was accounted for by applying a Bonferroni
317 correction;⁴⁶ the individual locus effect of a variant (kmer or PC) was considered significant if
318 its P value was smaller than α/n_p , where we took $\alpha = 0.05$ to be the genome-wide false-positive
319 rate and n_p to be the number of kmers or PCs with unique phylogenetic patterns, that is, unique
320 partitions of individuals according to allele membership. We identified 236627 unique kmer
321 patterns and 518 PCs, giving thresholds of 2.1×10^{-7} and 9.7×10^{-5} respectively.

322 **Data availability.** Sequence data has been submitted to Short Read Archive (Bioproject ID
323 PRJNA418899).

324
325

326 **Acknowledgments:** The authors would like to thank study participants. This study was funded
327 by the Wellcome Trust (MORU Grants 089275/H/09/Z and 089275/Z/09/Z), and a University
328 of Oxford Medical Research Fund awarded to C.E.M. (MRF/MT2015/2180). D.J.W. is a Sir
329 Henry Dale Fellow, jointly funded by the Wellcome Trust and the Royal Society (Grant
330 101237/Z/13/Z). B.C.Y. is a Research Training Fellow funded by the Wellcome Trust (Grant
331 101611/Z/13/Z). D.H.W. was funded by the National Institute for Health Research (NIHR)
332 Oxford Biomedical Research Centre (BRC) and the European Union's Seventh Framework
333 Programme under the grant agreement number 601783 (BELLEROPHON project). N.S. is
334 funded by a Public Health England (PHE)/University of Oxford Clinical Lectureship. K.E. was
335 funded by an academic clinical fellowship which was provided by the UK NIHR through the
336 University of Oxford. This research was supported by Core funding to the Wellcome Centre for
337 Human Genetics provided by the Wellcome Trust (090532/Z/09/Z). The views expressed are
338 those of the author(s) and not necessarily those of the NHS, PHE, the NIHR or the Department
339 of Health.

340
341 **Author contributions:** NPJD, CMP and CEM designed the study. SS, PS, VK, SH, VS, RB,
342 NS, KE, CMP, EN, PT & CEM collected bacterial samples and clinical data. CEM performed
343 DNA extraction for Illumina sequencing. LB performed Nanopore sequencing. BCY, SGE and
344 NDS performed bioinformatics on the study. BCY, SGE, DW, NPJD, DJW and CEM analysed
345 the data. RB and DC assisted with interpretation of findings. BCY, DJW and CEM wrote the
346 manuscript.

347
348 **Competing interests:** None to declare

349
350

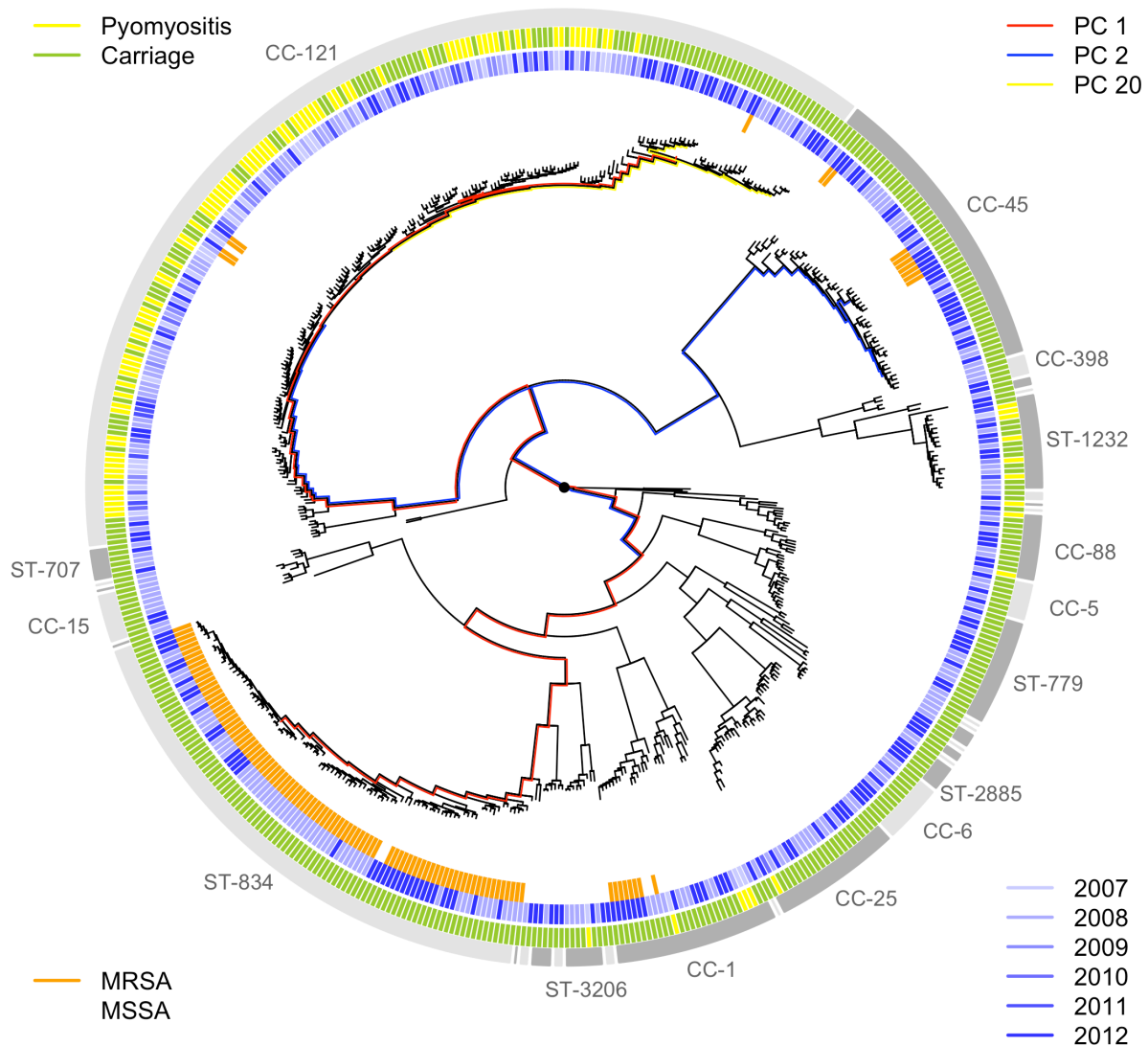
351 **References:**

- 352 1. Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ, Bentley SD, Maiden MC,
353 Parkhill J, Falush D. Genome-wide association study identifies vitamin B5 biosynthesis as a host
354 specificity factor in *Campylobacter*. *Proc Natl Acad Sci U S A*. **110(29)**, 11923-7 (2013) doi:
355 10.1073/pnas.1305559110.
- 356 2. Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, Mather AE, Hanage WP,
357 Goldblatt D, Nosten FH, Turner C, Turner P, Bentley SD, Parkhill J. Comprehensive identification
358 of single nucleotide polymorphisms associated with beta-lactam resistance within pneumococcal
359 mosaic genes. *PLoS Genet*. **10(8)**, e1004547 (2014) doi: 10.1371/journal.pgen.1004547.
- 360 3. Earle SG, Wu CH, Charlesworth J, Stoesser N, Gordon NC, Walker TM, Spencer CCA, Iqbal Z,
361 Clifton DA, Hopkins KL, Woodford N, Smith EG, Ismail N, Llewelyn MJ, Peto TE, Crook DW,
362 McVean G, Walker AS, Wilson DJ. Identifying lineage effects when controlling for population
363 structure improves power in bacterial association studies. *Nat Microbiol*. **1**, 16041 (2016) doi:
364 10.1038/nmicrobiol.2016.41.
- 365 4. Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, Croucher NJ, Marttinen P, Davies
366 MR, Steer AC, Tong SY, Honkela A, Parkhill J, Bentley SD, Corander J. Sequence element
367 enrichment analysis to determine the genetic basis of bacterial phenotypes. *Nat Commun*. **7**, 12797
368 (2016) doi: 10.1038/ncomms12797.
- 369 5. Falush D. Bacterial genomics: Microbial GWAS coming of age. *Nat Microbiol*. **1**, 16059 (2016) doi:
370 10.1038/nmicrobiol.2016.59.
- 371 6. Power R, Parkhill J, de Oliveira T. Microbial Genome-Wide Association Studies: Lessons from
372 Human GWAS, *Nat Rev Genet*. **18**, 41-50 (2017) doi: 10.1038/nrg.2016.132.
- 373 7. Chauhan S, Jain S, Varma S, Chauhan SS. Tropical pyomyositis (myositis tropicans): current
374 perspective. *Postgrad Med J*. **80(943)**, 267-70. (2004).
- 375 8. Verma S. Pyomyositis in Children. *Curr Infect Dis Rep*. **18(4)**, 12 (2016) doi: 10.1007/s11908-016-
376 0520-2.
- 377 9. Bickels J, Ben-Sira L, Kessler A, Wientroub S. Primary pyomyositis. *J Bone Joint Surg Am*. **84-**
378 **A(12)**, 2277-86 (2002).
- 379 10. Moriarty P, Leung C, Walsh M, Nourse C. Increasing pyomyositis presentations among children in
380 Queensland, Australia. *Pediatr Infect Dis J*. **34(1)**, 1-4 (2015) doi: 10.1097/INF.0000000000000470.
- 381 11. Pannaraj PS, Hulten KG, Gonzalez BE, Mason EO Jr, Kaplan SL. Infective pyomyositis and
382 myositis in children in the era of community-acquired, methicillin-resistant *Staphylococcus aureus*
383 infection. *Clin Infect Dis*. **43(8)**, 953-60 (2006).
- 384 12. Borges AH, Faragher B, Lalloo DG. Pyomyositis in the upper Negro river basin, Brazilian
385 Amazonia. *Trans R Soc Trop Med Hyg*. **106(9)**, 532-7 (2012) doi: 10.1016/j.trstmh.2012.06.008.
- 386 13. Shallcross LJ, Fragaszy E, Johnson AM, Hayward AC. The role of the Pantone-Valentine leucocidin
387 toxin in staphylococcal disease: a systematic review and meta-analysis. *Lancet Infect Dis*. **13**, 43-54
388 (2013) doi: 10.1016/S1473-3099(12)70238-4.
- 389 14. Vandenesch F, Naimi T, Enright MC, Lina G, Nimmo GR, Heffernan H, Liassine N, Bes M,
390 Greenland T, Reverdy ME, Etienne J. Community-acquired methicillin-resistant *Staphylococcus*
391 *aureus* carrying Pantone-Valentine leucocidin genes: worldwide emergence. *Emerg Infect Dis*. **9(8)**,
392 978-84 (2003).
- 393 15. Labandeira-Rey M, Couzon F, Boisset S, Brown EL, Bes M, Benito Y, Barbu EM, Vazquez V,
394 Höök M, Etienne J, Vandenesch F, Bowden MG. *Staphylococcus aureus* Pantone-Valentine
395 leucocidin causes necrotizing pneumonia. *Science*. **315(5815)**, 1130-3 (2007).
- 396 16. Villaruz AE, Bubeck Wardenburg J, Khan BA, Whitney AR, Sturdevant DE, Gardner DJ, DeLeo
397 FR, Otto M. A point mutation in the agr locus rather than expression of the Pantone-Valentine

- 398 leukocidin caused previously reported phenotypes in *Staphylococcus aureus* pneumonia and gene
399 regulation. *J Infect Dis.* **200(5)**, 724-34 (2009) doi: 10.1086/604728.
- 400 17. Löffler B, Hussain M, Grundmeier M, Brück M, Holzinger D, Varga G, Roth J, Kahl BC, Proctor
401 RA, Peters G. *Staphylococcus aureus* Panton-Valentine leukocidin is a very potent cytotoxic factor
402 for human neutrophils. *PLoS Pathog.* **6(1)**, e1000715 (2010) doi: 10.1371/journal.ppat.1000715.
- 403 18. Boakes E, Kearns AM, Ganner M, Perry C, Hill RL, Ellington MJ. Distinct bacteriophages encoding
404 Panton-Valentine leukocidin (PVL) among international methicillin-resistant *Staphylococcus aureus*
405 clones harboring PVL. *J Clin Microbiol.* **49(2)**, 684-92 (2011) doi: 10.1128/JCM.01917-10.
- 406 19. McCarthy AJ, Witney AA, Lindsay JA. *Staphylococcus aureus* temperate bacteriophage: carriage
407 and horizontal gene transfer is lineage associated. *Front Cell Infect Microbiol.* **2**, 6 (2012) doi:
408 10.3389/fcimb.2012.00006.
- 409 20. Sina H, Ahoyo TA, Moussaoui W, Keller D, Bankolé HS, Barogui Y, Stienstra Y, Kotchoni SO,
410 Prévost G, Baba-Moussa L. Variability of antibiotic susceptibility and toxin production of
411 *Staphylococcus aureus* strains isolated from skin, soft tissue, and bone related infections. *BMC*
412 *Microbiol.* **13**, 188 (2013) doi:10.1186/1471-2180-13-188.
- 413 21. García C, Hallin M, Deplano A, Denis O, Sihuíncha M, de Groot R, Gotuzzo E, Jacobs J.
414 *Staphylococcus aureus* causing tropical pyomyositis, Amazon Basin, Peru. *Emerg Infect Dis.* **19(1)**,
415 123-5 (2013) doi: 10.3201/eid1901.120819.
- 416 22. Otto M. A MRSA-terious enemy among us: end of the PVL controversy? *Nat Med.* **17(2)**, 169-70
417 (2011) doi: 10.1038/nm0211-169.
- 418 23. Day NPJ. Panton-Valentine leucocidin and staphylococcal disease *Lancet Infect Dis.* **13**: 5-6 (2013)
419 doi: 10.1016/S1473-3099(12)70265-7.
- 420 24. Zhou X and Stephens M. Genome-wide efficient mixed-model analysis for association studies.
421 *Nature Genetics.* **44**, 821–824 (2012) doi: 10.1038/ng.2310.
- 422 25. Diep BA, Palazzolo-Ballance AM, Tattévin P, Basuino L, Braughton KR, Whitney AR, Chen L,
423 Kreiswirth BN, Otto M, DeLeo FR, Chambers HF. Contribution of Panton-Valentine leukocidin in
424 community-associated methicillin-resistant *Staphylococcus aureus* pathogenesis. *PLoS One.* **3(9)**,
425 e3198 (2008) doi: 10.1371/journal.pone.0003198.
- 426 26. Bocchini CE, Hultén KG, Mason EO Jr, Gonzalez BE, Hammerman WA, Kaplan SL. Panton-
427 Valentine leukocidin genes are associated with enhanced inflammatory response and local disease in
428 acute hematogenous *Staphylococcus aureus* osteomyelitis in children. *Pediatrics.* **117(2)**, 433-40
429 (2006).
- 430 27. Kurt K, Rasigade JP, Laurent F, Goering RV, Žemličková H, Machova I, Struelens MJ, Zautner AE,
431 Holtfreter S, Bröker B, Ritchie S, Reaksmey S, Limmathurotsakul D, Peacock SJ, Cuny C, Layer F,
432 Witte W, Nübel U. Subpopulations of *Staphylococcus aureus* Clonal Complex 121 Are Associated
433 with Distinct Clinical Entities. *PLoS ONE.* **8(3)**, e58155 (2013) doi: 10.1371/journal.pone.0058155.
- 434 28. Nickerson EK, Wuthiekanun V, Kumar V, Amornchai P, Wongdeethai N, Chheng K, Chantratita N,
435 Putschhat H, Thaipadungpanit J, Day NP, Peacock SJ. Emergence of community-associated
436 methicillin-resistant *Staphylococcus aureus* carriage in children in Cambodia. *Am J Trop Med Hyg.*
437 **84(2)**, 313-7 (2011) doi: 10.4269/ajtmh.2011.10-0300.
- 438 29. Landrum ML, Lalani T, Niknian M, Maguire JD, Hospenthal DR, Fattom A, Taylor K, Fraser J,
439 Wilkins K, Ellis MW, Kessler PD, Fahim RE, Tribble DR. Safety and immunogenicity of a
440 recombinant *Staphylococcus aureus* α -toxoid and a recombinant Panton-Valentine leukocidin
441 subunit, in healthy adults. *Hum Vaccin Immunother.* **13(4)**, 791-801 (2017) doi:
442 10.1080/21645515.2016.1248326.
- 443 30. Saeed K, Gould I, Esposito S, Ahmad-Saeed N, Ahmed SS, Alp E, Bal AM, Bassetti M, Bonnet E,
444 Chan M, Coombs G, Dancer SJ, David MZ, De Simone G, Dryden M, Guardabassi L, Hanitsch LG,
445 Hijazi K, Krüger R, Lee A, Leistner R, Pagliano P, Righi E, Schneider-Burrus S, Skov RL, Tattévin

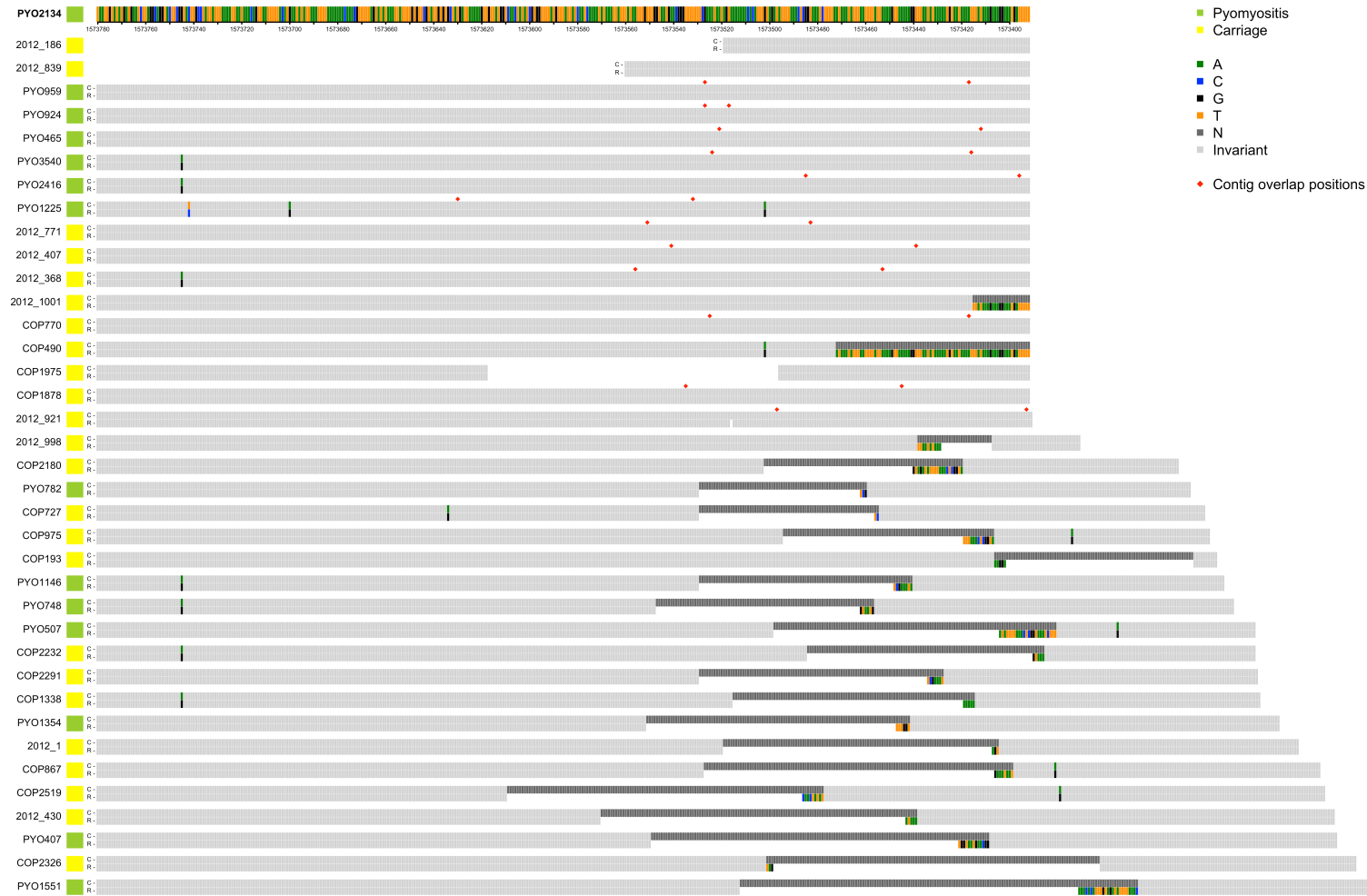
- 446 P, Van Wamel W, Vos MC, Voss A; International Society of Chemotherapy. Pantone-Valentine
447 leukocidin-positive *Staphylococcus aureus*: a position statement from the International Society of
448 Chemotherapy. *Int J Antimicrob Agents*. 2018 Jan;51(1):16-25. doi:
449 10.1016/j.ijantimicag.2017.11.002.
- 450 31. CLSI. *Performance Standards for Antimicrobial Susceptibility Testing*. 24th ed. CLSI supplement
451 M100. Wayne, PA: Clinical and Laboratory Standards Institute; 2014.
- 452 32. Zerbino DR and Birney E. Velvet: Algorithms for de novo short read assembly using de bruijn
453 graphs *Genome Res*. **18(5)**, 821-9 (2008) doi: 10.1101/gr.074492.107.
- 454 33. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool *J Mol*
455 *Biol*. **215(3)**, 403-10 (1990).
- 456 34. Gordon NC, Price JR, Cole K, Everitt R, Morgan M, Finney J, Kearns AM, Pichon B, Young B,
457 Wilson DJ, Llewelyn MJ, Paul J, Peto TE, Crook DW, Walker AS, Golubchik T. Prediction of
458 *Staphylococcus aureus* antimicrobial resistance by whole-genome sequencing *J Clin Microbiol*.
459 **52(4)**,1182-91 (2014) doi: 10.1128/JCM.03117-13.
- 460 35. Lunter G and Goodson M. Stampy: A statistical algorithm for sensitive and fast mapping of illumina
461 sequence reads *Genome Res*. **21(6)**, 936-9 (2011) doi: 10.1101/gr.111120.110.
- 462 36. Diep BA, Gill SR, Chang RF, Phan TH, Chen JH, Davidson MG, Lin F, Lin J, Carleton HA,
463 Mongodin EF, Sensabaugh GF, Perdreau-Remington F. Complete genome sequence of USA300, an
464 epidemic clone of community-acquired methicillin-resistant *Staphylococcus aureus*. *Lancet*.
465 **367(9512)**, 731-9 (2006).
- 466 37. Didelot X, Eyre DW, Cule M, Ip CL, Ansari MA, Griffiths D, Vaughan A, O'Connor L, Golubchik
467 T, Batty EM, Piazza P, Wilson DJ, Bowden R, Donnelly PJ, Dingle KE, Wilcox M, Walker AS,
468 Crook DW, Peto TE, Harding RM. Microevolutionary analysis of *Clostridium difficile* genomes to
469 investigate transmission. *Genome Biol*. **13(12)**,R118 (2012) doi: 10.1186/gb-2012-13-12-r118.
- 470 38. Young BC, Golubchik T, Batty EM, Fung R, Larner-Svensson H, Votintseva AA, Miller RR,
471 Godwin H, Knox K, Everitt RG, Iqbal Z, Rimmer AJ, Cule M, Ip CL, Didelot X, Harding RM,
472 Donnelly P, Peto TE, Crook DW, Bowden R, Wilson DJ. Evolutionary dynamics of *Staphylococcus*
473 *aureus* during progression from carriage to disease. *Proc Natl Acad Sci U S A*. **109(12)**, 4550 (2012)
474 doi: 10.1073/pnas.1113219109.
- 475 39. Golubchik T, Batty EM, Miller RR, Farr H, Young BC, Larner-Svensson H, Fung R, Godwin H,
476 Knox K, Votintseva A, Everitt RG, Street T, Cule M, Ip CL, Didelot X, Peto TE, Harding RM,
477 Wilson DJ, Crook DW, Bowden R. Within- host evolution of *Staphylococcus aureus* during
478 asymptomatic carriage. *PLoS One*. **8(5)**: e61319 (2013) doi: 10.1371/journal.pone.0061319.
- 479 40. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies from
480 short and long sequencing reads. *PLoS Comput Biol*. **13(6)**:e1005595. (2017) doi:
481 10.1371/journal.pcbi.1005595.
- 482 41. Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large
483 phylogenies *Bioinformatics*. **30(9)**,1312-3 (2014) doi: 10.1093/bioinformatics/btu033.
- 484 42. Didelot X and Wilson DJ. ClonalFrameML: Efficient inference of recombination in whole bacterial
485 genomes *PLoS Comput Biol*. **11(2)**, e1004041 (2015) doi: 10.1371/journal.pcbi.1004041.
- 486 43. Rizk G, Lavenier D, and Chikhi R. DSK: k-mer counting with very low memory usage.
487 *Bioinformatics*. **29**, 652–653 (2013) doi: 10.1093/bioinformatics/btt020.
- 488 44. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. *Nat Methods*. **4**; 9(4), 357-9
489 (2012) doi: 10.1038/nmeth.1923.
- 490 45. Darling AC, Mau B, Blattner FR, Perna NT. Mauve: Multiple alignment of conserved genomic
491 sequence with rearrangements. *Genome Res*. **14(7)**, 1394-403 (2004) doi: 10.1101/gr.2289704
- 492 46. Dunn OJ. Estimation of the medians for dependent variables. *Ann. Math. Stat*. **30**, 192–197 (1959).

493 **Supplementary information**

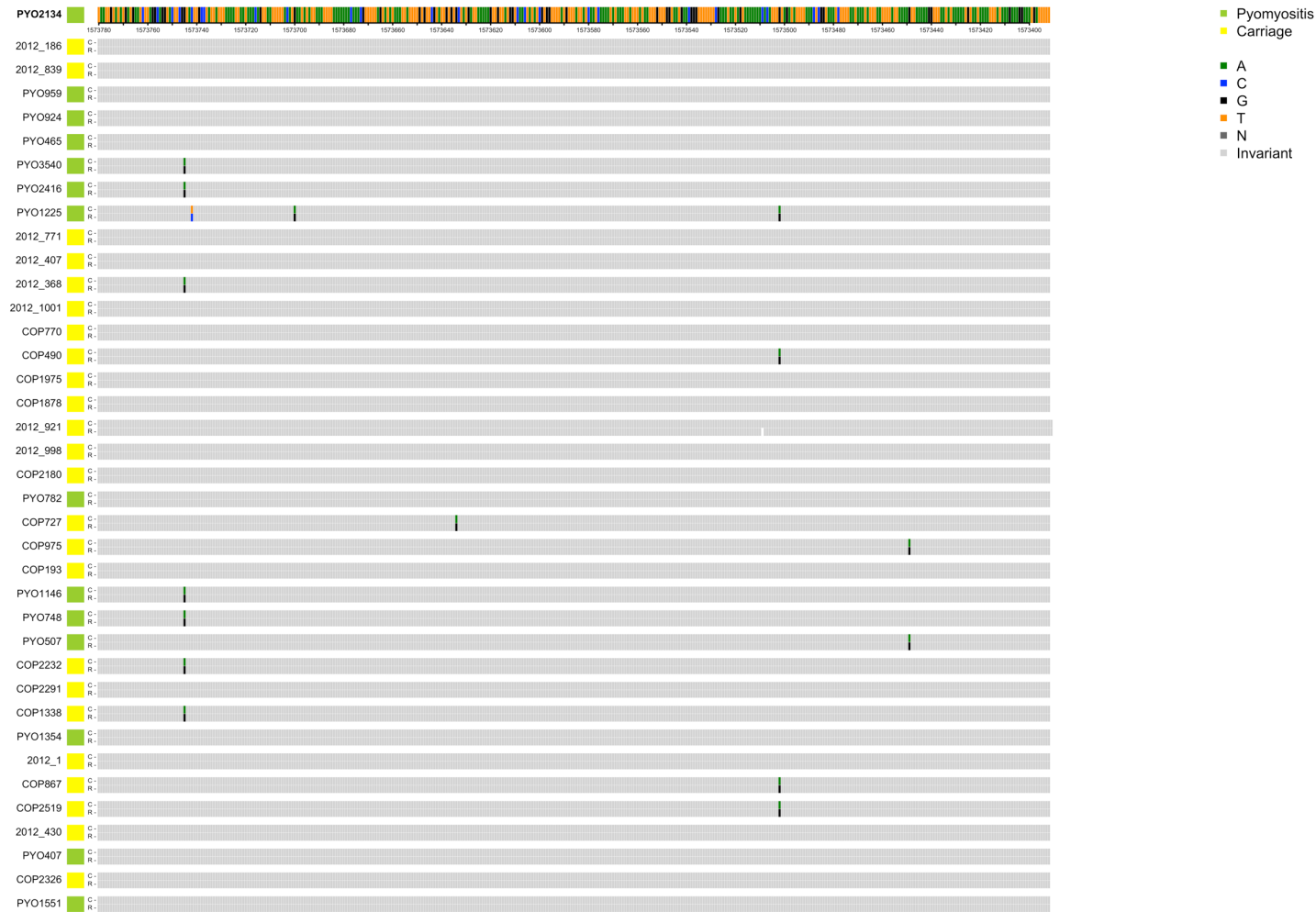


494

495 **Figure S1** Sampling frequencies of the major strains were stable over time. The year of
496 sampling (2007-2008, blue shaded lines) and MRSA status (orange lines) are illustrated around
497 the phylogeny of the bacteria sampled from pyomyositis cases (gold lines) and asymptomatic
498 carriage controls (green lines). The three PCs most significantly associated with case/control
499 status are also shown (PCs 1, 2 and 20 by red, blue and yellow branches respectively).



500
 501 **Figure S2** Alignments of reference genome PYO2134 assembly (R) with 37 *de novo* assemblies of Illumina short-read sequencing (C) which
 502 feature either ambiguities (Ns) or contig boundaries in the region 389bp upstream of PVL coding sequence. Contig boundaries, when
 503 overlapping, are marked with a red diamond. Ns in the assembly are marked in dark grey. Polymorphisms are colour-coded by base.



504
 505 **Figure S3** Alignments of reference genome PYO2134 assembly (R) with 37 *de novo* hybrid assemblies combining Oxford Nanopore long-read
 506 and Illumina short-read sequencing (C) which featured ambiguities (Ns) or hybrid assembly contig boundaries in the region 389bp upstream of
 507 PVL coding sequence in the Illumina short-read only assemblies. All ambiguities and contig boundaries are resolved in the hybrid assemblies.

508

| | Pyomyositis (2007-2012) | Nasal carriage (2008) | Nasal carriage (2012) |
|--------------------|----------------------------|--------------------------|--------------------------|
| Number of isolates | 101 | 222 | 195 |
| Age (med, IQR) | 7.8 years (4.2-11.8) | 5.9 years (2.5- 8.3) | 6.3 years (4.1- 9.9) |
| Male (n (%)) | 66/97 (68%) | 122/221 (55.2%) | 105/195 (53.8%) |
| MRSA (n (%)) | 0 (0%) | 61 (27.5%) | 52 (26.7%) |

509 **Table S1.** Isolates included in this study.

510 **Table S2:** All significant kmers from short read sequencing assembly, evidence of
511 association, frequency and best match on blastn to all *S. aureus* sequences in Genbank.