

1 **Panton-Valentine leucocidin is the key determinant of *Staphylococcus aureus***  
2 **pyomyositis in a bacterial GWAS**

3 Bernadette C Young<sup>1,2</sup>, Sarah G Earle<sup>1</sup>, Sona Soeng<sup>3</sup>, Poda Sar<sup>3</sup>, Varun Kumar<sup>4</sup>, Songly Hor<sup>3</sup>,  
4 Vuthy Sar<sup>3</sup>, Rachel Bousfield<sup>5</sup>, Nicholas D Sanderson<sup>1</sup>, Leanne Barker<sup>1</sup>, Nicole Stoesser<sup>1,6</sup>,  
5 Katherine RW Emary<sup>2</sup>, Christopher M Parry<sup>7,8</sup>, Emma K Nickerson<sup>5</sup>, Paul Turner<sup>3,9</sup>, Rory  
6 Bowden<sup>10</sup>, Derrick Crook<sup>1,2,6</sup>, David Wyllie<sup>1,6,11</sup>, Nicholas PJ Day<sup>9,13</sup>†, Daniel J Wilson<sup>1,10,12</sup> †,  
7 Catrin E Moore<sup>9,13</sup>†\*

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

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

12 <sup>3</sup> Cambodia Oxford Medical Research Unit, Angkor Hospital for Children, Siem Reap,  
13 Cambodia

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

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

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

20 <sup>7</sup> Clinical Sciences, Liverpool School of Tropical Medicine, Liverpool, L3 5QA

21 <sup>8</sup> School of Tropical Medicine and Global Health, Nagasaki University, Nagasaki, Japan

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

24 <sup>10</sup> Wellcome Centre for Human Genetics, University of Oxford, Oxford, OX3 7BN, UK

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

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

29 <sup>13</sup> 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](mailto: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, predominantly caused by *Staphylococcus aureus*. To understand  
37 the contribution of bacterial genomic factors to pyomyositis, we conducted a genome-wide  
38 association study of *S. aureus* cultured from 101 children with pyomyositis and 417 children  
39 with asymptomatic nasal carriage attending the Angkor Hospital for Children, Cambodia. We  
40 found a strong relationship between bacterial genetic variation and pyomyositis, with estimated  
41 heritability 63.8% (95% CI 49.2-78.4%). The presence of the Pantone-Valentine leucocidin  
42 (PVL) locus increased the odds of pyomyositis 130-fold ( $p=10^{-17.9}$ ). The signal of association  
43 mapped both to the PVL-coding sequence and the sequence immediately upstream. Together  
44 these regions explained over 99.9% of heritability (95% CI 93.5-100%). Our results establish  
45 staphylococcal pyomyositis, like tetanus and diphtheria, as critically dependent on a single  
46 toxin and demonstrate the potential for association studies to identify specific bacterial genes  
47 promoting severe human disease.  
48

## 49 Introduction

50 Microbial genome sequencing and bacterial genome-wide association studies (GWAS) present  
51 new opportunities to discover bacterial genes involved in the pathogenesis of serious  
52 infections.<sup>1-6</sup> Pyomyositis is a severe infection of skeletal muscle most commonly seen in  
53 children in the tropics.<sup>7-9</sup> In up to 90% of cases it is caused by a single bacterial pathogen,  
54 *Staphylococcus aureus* (*S. aureus*).<sup>7-10</sup> Unlike infections of the skin and superficial soft tissues,  
55 the skin and subcutaneous tissues are not usually involved in pyomyositis, by contrast to  
56 intense inflammation in the infected muscles.<sup>7,8</sup> Pyomyositis is thought to arise from  
57 haematogenous seeding of bacteria from blood to muscle.<sup>8</sup> There is evidence that some  
58 *S. aureus* strains have heightened propensity to cause pyomyositis – the incidence in the USA  
59 doubled during an epidemic of community-associated methicillin resistant *S. aureus* (CA-  
60 MRSA)<sup>11</sup> – but molecular genetic investigation of *S. aureus* from pyomyositis has been  
61 limited.<sup>12</sup>

62 Panton-Valentine leucocidin (PVL), a well-known staphylococcal toxin causing purulent skin  
63 infections and found in epidemics caused by CA-MRSA, has been implicated in pyomyositis,  
64 pneumonia and other *S. aureus* disease manifestations, but its role in these invasive infections  
65 is disputed.<sup>13-16</sup> PVL is a bipartite pore-forming toxin comprising the co-expressed LukF-PV  
66 and LukS-PV proteins.<sup>17,18</sup> The coding sequence for PVL, *lukSF-PV*, is usually carried on  
67 bacteriophages,<sup>13,17</sup> which facilitate *lukSF-PV* exchange between lineages.<sup>19</sup> The mechanism of  
68 PVL toxicity has been shown to involve cell lysis in human myeloid cells, particularly  
69 neutrophils, by insertion into the cellular membrane,<sup>21</sup> leading the tissue to release  
70 inflammatory neutrophil products.<sup>22</sup> Neutrophil lysis is mediated by PVL binding to target  
71 complement receptors C5aR; in binding, PVL has both toxic and immunomodulatory effects, as  
72 it also inhibits C5a mediated immune activation.<sup>23</sup>

73 Although small case series testing for candidate genes have reported a high prevalence of PVL  
74 among pyomyositis-causing *S. aureus*,<sup>11,23,24</sup> a detailed meta-analysis found no evidence for an  
75 increased rate of musculoskeletal infection (or other invasive disease) in PVL-positive bacteria  
76 versus controls.<sup>13</sup> These conflicting results may reflect insufficiently powered studies, and  
77 some case series lack comparative control strains.<sup>23</sup> A further problem with the use of candidate  
78 gene studies in studying pathogenesis is that they may miss important variation elsewhere in  
79 the genome. One study reporting a critical role for PVL in the causation of severe pneumonia<sup>15</sup>  
80 was later found to have overlooked mutations in key regulatory genes, capable of producing the  
81 virulence that had been attributed to PVL by the original study.<sup>16</sup> Thus, while some evidence  
82 suggests an association between pyomyositis and PVL, there remains significant uncertainty  
83 regarding the bacterial genetic predisposition of *S. aureus* to pyomyositis, and whether PVL is  
84 an important virulence factor, or merely an epiphenomenon, carried by bacteria alongside  
85 unidentified genetic determinants.<sup>25,26</sup>

86 GWAS offer a means to screen entire bacterial genomes to discover genes and genetic variants  
87 associated with disease risk. They are particularly appealing because they enable the  
88 investigation of traits not readily studied in the laboratory, and do not require the nomination of  
89 specific candidate genes.<sup>5</sup> Proof-of-principle GWAS in bacteria have demonstrated the  
90 successful rediscovery of known antimicrobial resistance (AMR) determinants.<sup>2,3,4</sup> However,  
91 AMR is under extraordinarily intense selection in bacteria. More subtle traits, including host  
92 specificity<sup>1</sup> and the duration of pneumococcal carriage,<sup>27</sup> have also been demonstrated using  
93 GWAS. Promising results for GWAS in human infecting bacteria include identifying possible  
94 loci for invasive infection with *Streptococcus pyogenes*<sup>4</sup> and *Staphylococcus epidermidis*<sup>28</sup>, and  
95 the identification of virulence associated genes corresponding with regional differences in  
96 disease manifestations of melioidosis.<sup>29</sup> Within species, lineage specific variants have been

97 shown to predict mortality following *S. aureus* bacteraemia.<sup>30</sup> These studies support the  
98 potential GWAS has to precisely pinpoint genes and genetic variants underlying the propensity  
99 to cause specific human infections, making it a promising tool to investigate the possible  
100 contribution of bacterial genomic variation to pyomyositis.

## 101 **Results**

102 To understand the bacterial genetic basis of pyomyositis, we sampled and whole-genome  
103 sequenced *S. aureus* from 101 pyomyositis infections and 417 asymptomatic nasal carriage  
104 episodes in 518 children attending Angkor Hospital for Children in Siem Reap, Cambodia  
105 between 2008 and 2012 (Supplementary File 1). As expected, we observed representatives of  
106 multiple globally common lineages in Cambodia, together with some globally less common  
107 lineages at high frequency, in particular clonal complex (CC) 121, identified by multi-locus  
108 sequence typing (MLST). There were no major changes in lineage frequency over time (Figure  
109 1-figure supplement 1).

110 In our study, some *S. aureus* lineages were strongly overrepresented among cases of  
111 pyomyositis compared with asymptomatic, nasally-carried controls over the same time period.  
112 Notably, 86/101 (85%) of pyomyositis cases were caused by CC-121 bacteria, whereas no  
113 pyomyositis cases were caused by the next two most commonly carried lineages, sequence type  
114 (ST)-834 and CC-45 (Figure 1). We estimated the overall heritability of case/control status to  
115 be 63.8% (95% CI 49.2-78.4%) in the sample, reflecting the strong relationship between  
116 bacterial genetic variation and case/control status. We used *bugwas*<sup>6</sup> to decompose this  
117 heritability into the principal components (PCs) of bacterial genetic variation. PC 1, which  
118 distinguished CC-121 (the most common pyomyositis lineage) from ST-834 (which was only  
119 found in carriage), showed the strongest association with case /control status ( $p = 10^{-29.6}$ , Wald  
120 test). The strongest associations were with PC 20, which differentiated a sub-lineage of CC-121  
121 within which no cases were seen ( $p = 10^{-13.9}$ ), and PC 2, which distinguished CC-45 from the  
122 rest of the species ( $p = 10^{-4.9}$ ).

123 We conducted a GWAS to identify bacterial genetic variants associated with pyomyositis,  
124 controlling for differences in pyomyositis prevalence between *S. aureus* lineages. We used a  
125 kmer-based approach<sup>1</sup> in which every variably present 31bp DNA sequence observed among  
126 the 518 genomes was tested for association with pyomyositis *versus* asymptomatic nasal  
127 carriage, controlling for population structure using GEMMA.<sup>31</sup> These kmers captured bacterial  
128 genetic variation including single nucleotide polymorphisms (SNPs), insertions or deletions  
129 (indels), and presence or absence of entire accessory genes. We found 10.7 million unique  
130 kmers variably present across the bacterial genomes. In total, 9,175 kmers were significantly  
131 associated with case/control status after correction for multiple testing ( $10^{-6.8} \leq p \leq 10^{-21.4}$ ;  
132 Figure 2A). When mapped to the *de novo* assembly of a CC121 isolate from pyomyositis  
133 (PYO2134), the vast majority of these kmers (9,074/9,175; 98.9%) localised to a 45.7kb region  
134 spanning an integrated prophage with 95% nucleotide identity to  $\phi$ SLT (Figure 2B). Most  
135 (9,173/9,175; 99.98%) significant kmers were found at an increased frequency in pyomyositis,  
136 with odds ratios (OR) ranging from 2.7 to 139.8, indicating that the presence of each was  
137 associated with increased risk of disease. The presence of bacteriophage  $\phi$ SLT was thus  
138 strongly associated with pyomyositis.

139 We were able to localise the most statistically significant signal of association to kmers that  
140 mapped within  $\phi$ SLT to the *lukS-PV* and *lukF-PV* cargo genes. These genes encode the  
141 subunits of PVL, which multimerise into a pore-forming toxin capable of rapidly lysing the  
142 membranes of human neutrophils.<sup>17,25</sup> 1,630 kmers tagging the presence of the *lukSF-PV*  
143 coding sequences (CDS) were highly significantly associated with disease, being present in  
144 98/101 (97%) pyomyositis cases and 84/417 (20%) carriage controls (unadjusted OR 129.5,

145  $p=10^{-17.9}$ ). Kmers tagging variation in the 389bp region immediately upstream of the CDS were  
146 also strongly associated with disease ( $p=10^{-21.4}$ ). The most significant of these kmers were co-  
147 present with the CDS in the same cases (98/101, 97.0%), but present in fewer controls (79/417,  
148 18.9%), producing an OR of 140.

149 Closer examination of this ~400bp upstream region in genomes assembled from short-read  
150 Illumina sequencing showed that assembly of the region was problematic, with breaks or gaps  
151 in the assembly (Figure 2-figure supplement 2). To improve the accuracy of this region of the  
152 assembled genomes we performed long-read Oxford Nanopore sequencing on the 37 genomes  
153 with incomplete or discontinuous assembly upstream of the PVL CDS. By integrating long-  
154 read and short-read data we were able to assemble a single contig spanning this region in all  
155 isolates (Figure 2-figure supplement 2). When these improved assemblies were introduced, the  
156 signal of association upstream of the PVL CDS was no more significant than within the CDS  
157 (Figure 2C). Therefore, the presence of genomic sequence spanning the PVL toxin-coding  
158 sequences and the upstream, presumed regulatory, region exhibited the strongest association  
159 with pyomyositis in the *S. aureus* genome. All isolates with kmers mapping to the PVL CDS  
160 had 98% or more coverage for the PVL CDS genes in *de novo* assembly (Figure 2-figure  
161 supplement 3). The signal of association in the earlier and later periods of the study were  
162 examined and found to be consistent (Figure 2-figure supplement 4).

163 Out of 9,175 kmers significantly associated with pyomyositis, we only found 101 kmers related  
164 to regions outside the PVL-carrying prophage (Figure 2A, Supplementary File 2). Two kmers  
165 mapping at a position near 0.2Mb in the PYO2014 reference genome showed homology to  
166 platelet adhesin *sraP* by BLAST. Thirty-five kmers mapping to a 50bp non-coding fragment at  
167 0.6Mb and two kmers mapping to 2.8Mb showed homology to an MSSA476 intergenic  
168 sequence between adhesin-encoding *sdrC* and *sdrD* by BLAST. One kmer mapping to position  
169 2.0Mb showed no sequence homology by BLAST. Sixty-one kmers mapping to a 61bp non-  
170 coding region at 2.7Mb showed homology to an MSSA476 intergenic sequence between  
171 acetyltransferase-encoding genes SAS2453 and SAS2454 by BLAST. In conclusion, these  
172 other signals were short, fragmentary and mostly non-coding so we did not investigate them  
173 further.

174 The presence of high-risk kmers mapping to the PVL region explained the vast majority of  
175 observed heritability. When the presence or absence of the most significant kmer pattern, a set  
176 of kmers with an identical pattern of presence in the population, all of which mapped to the  
177 PVL region, was included as a covariate in GEMMA, the remaining heritability not explained  
178 by other factors was estimated and found to be 0.0% (95% CI 0-2.5%). Thus, the point estimate  
179 for heritability (not explained by the inclusion of PVL-tagging kmers) is reduced by 100%  
180 (95% CI 93.5-100%), meaning we have little evidence for any remaining heritability in  
181 case/control status.

182 Presence or absence of the PVL region accounted for the differences in pyomyositis rates  
183 between lineages. It was common in pyomyositis-associated lineages including CC-121 and  
184 absent from carriage/non-pyomyositis-associated lineages including ST-834 and CC-45  
185 (Figure. 1), explaining over 99.9% of observed heritability in case-control status. It was  
186 infrequent in the non-pyomyositis-associated sub-lineage of CC-121 (2/36, 5.6%), and  
187 sporadically present in pyomyositis cases in otherwise non-pyomyositis-associated, PVL-  
188 negative strains CC-1 and CC-88. Its absence from only three cases (in lineages CC-88, CC-1  
189 and CC-121) suggested that the PVL region approached necessity for development of  
190 pyomyositis in the current setting in Cambodian children, while its presence in 20% of controls  
191 indicated that PVL-associated pyomyositis is incompletely penetrant, i.e. presence of the PVL  
192 region does not always lead to disease.

193 PVL genes were carried on multiple genetic backgrounds in this population, and the phage  
194 backgrounds vary by clonal complex. We examined all assemblies for sequence similarity to  
195 six known bacteriophages that carry the PVL genes,<sup>18</sup> as well as the 45.7kb region in  
196 PYO2134, a hypothesised integrated prophage identified in the reference genome prepared for  
197 this study, which we have called  $\phi$ CC121 (Figure S4). The finding of PVL genes on BLAST  
198 corresponded completely with the presence of kmers mapping across the PVL locus. We find  
199 sequences with >95% homology to four of these seven phages in the population. Regions in  
200 some assemblies showed homology for multiple phages, reflecting the similarity between  
201  $\phi$ SLT,  $\phi$ Sa2USA and  $\phi$ CC121 rather than the presence of multiple phages, and resolution of  
202 phages was limited by fragmented assemblies from short reads (Supplementary File 3). Phage  
203 types were restricted within most lineages, with  $\phi$ SLT found in ST-3206,  $\phi$ Sa2USA in ST-  
204 1232 and  $\phi$ PVL in CC-1.  $\phi$ CC121 was the phage most often identified in the dominant  
205 pyomyositis strain CC-121, but it was absent in all but one isolate from the low risk subclade  
206 within CC-121. Strikingly, PVL negative isolates in CC-121 and ST-834 strains frequently  
207 retained sequence homologous to >95% of  $\phi$ SLT, suggesting that some PVL-negative CC-121  
208 isolates lost PVL secondarily by gene deletion rather than prophage excision.

## 209 Discussion

210 In this study we found a strong association between pyomyositis, a highly distinctive tropical  
211 infection of skeletal muscle in children, and Panton-Valentine leukocidin, a bacterial toxin  
212 commonly carried by bacteriophages. We found that a single coding region together with the  
213 upstream sequence are all but necessary for the development of pyomyositis: its sporadic  
214 presence is associated with pyomyositis in otherwise low-frequency strains, and its absence is  
215 associated with asymptomatic carriage in a high-propensity strain. PVL appears to be carried  
216 on multiple phage backgrounds in this population, but PVL-positive lineages generally carry a  
217 single phage type, as expected given the observed strain restriction of phages in *S. aureus*.<sup>33,34</sup>  
218 The locally common PVL-positive CC-121 lineage contributes most strongly to the prevalence  
219 of pyomyositis in Cambodian children.

220 While PVL has long been thought an important *S. aureus* virulence factor,<sup>35-37</sup> its role in  
221 invasive disease has been controversial,<sup>25,26</sup> with conflicting results in case-control studies and  
222 an absence of supporting evidence on meta-analysis.<sup>13</sup> In previous studies the PVL positive  
223 USA300 lineage was associated with musculoskeletal infection (both pyomyositis and  
224 osteomyelitis), however in these studies almost all such infections were caused by the USA300  
225 strain, so the role of PVL was almost completely confounded by both methicillin resistance and  
226 strain background.<sup>11,36</sup> In our study, this confounding influence is broken down by the  
227 movement of PVL on mobile genetic elements (MGEs). Despite the emergence of CA-MRSA  
228 in carriage in Cambodia,<sup>38</sup> all the pyomyositis cases were MSSA (Figure 1-figure supplement  
229 1). By applying GWAS methods to a well-powered cohort, our study resolves the controversy  
230 around pyomyositis and PVL, demonstrating strong heritability which localises to a single  
231 region, even when the full bacterial genome is considered. Bacterial GWAS can pinpoint  
232 virulence variants when MGEs act to unravel linkage disequilibrium, if effect sizes are  
233 sufficiently strong.

234 There is strong biological plausibility for the association demonstrated in this study. PVL is a  
235 well characterised *S. aureus* toxin, toxic to myeloid cells, which form a first line of defence  
236 against bacterial infection,<sup>20</sup> and in binding to myeloid cells by a complement receptor (C5aR),  
237 exerts immunomodulatory effects.<sup>22</sup> The establishment of muscle abscesses is a critical step in  
238 the pathogenesis of pyomyositis, but unlike renal, hepatic and splenic abscesses, skeletal  
239 muscle abscesses are rarely observed in experimental models of bacteraemia, unless there is  
240 preceding muscle trauma.<sup>39</sup> *S. aureus* strains containing PVL show increased duration of

241 bacteraemia in a rabbit model of sepsis,<sup>33</sup> and result in larger muscle abscesses.<sup>40</sup> PVL has been  
242 found strongly bound to necrotic muscle in an individual with myositis associated with  
243 necrotizing fasciitis.<sup>41</sup> These observations support the hypothesis PVL may facilitate bacterial  
244 seeding to muscles via the bloodstream, and tropism for muscular infection.

245 These results establish that, for children in Cambodia, staphylococcal pyomyositis is a disease  
246 whose pathogenesis depends crucially on a single toxin. This property is shared by toxin-  
247 driven, vaccine-preventable diseases such as tetanus and diphtheria. Therefore, vaccines that  
248 generate neutralising anti-toxin antibodies against PVL<sup>42</sup> may protect human populations  
249 specifically against this common tropical disease. These results also raise the hypothesis that  
250 antibiotics which decrease toxin expression, and have been recommended in some PVL-  
251 associated infections,<sup>43</sup> may offer specific clinical benefit in treating pyomyositis. More  
252 generally, our study provides an example of how microbial GWAS can be used to elucidate the  
253 pathogenesis of bacterial infections and identify specific virulence genes associated with  
254 human disease.

255 **Materials and Methods:**

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

264 We identified *S. aureus* nasal colonisation from two cohort studies undertaken at Angkor  
265 Hospital for Children. The first were selected from a collection characterising nasal  
266 colonization in the region between September and October, 2008, which has previously been  
267 described using multi-locus sequence typing.<sup>38</sup> The swabs had been saved at  $-80^{\circ}\text{C}$  since the  
268 study, these samples were re-examined for the presence of *S. aureus* using selective agar,  
269 confirmed using Staphaurex (Remel, Lenexa, USA) and the DNase agar test (Oxoid,  
270 Hampshire, UK). Antimicrobial susceptibility testing was performed according to the 2014  
271 Clinical and Laboratory Standards Institute guidelines (M100-24).<sup>44</sup>

272 We undertook a second cohort study in 2012. Inclusion criteria were children ( $\leq 16$  years)  
273 attending as an outpatient at Angkor Hospital for Children with informed consent. There were  
274 no exclusion criteria. Children were swabbed between the 2-7th July 2012, using sterile cotton  
275 tipped swabs pre-moistened (with phosphate buffered saline, PBS) using 3 full rotations of the  
276 swab within the anterior portion of each nostril with one swab being used for both nostrils, the  
277 ends were broken into bottles containing sterile PBS and kept cool until plated in the laboratory  
278 (within the hour). The swabs were plated onto Mannitol Salt agar to select for *S. aureus*. The  
279 M100-24 CLSI<sup>44</sup> standards were followed for susceptibility testing and bacteria stored in  
280 tryptone soya broth and glycerol at  $-80^{\circ}\text{C}$ .

281 We selected controls from carriers in these two cohorts using the excel randomization function:  
282 222 of 519 from the 2008 cohort and 195 of 261 from the 2012 cohort.

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

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

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

299 **Variant calling.** For short-read sequencing we used Velvet<sup>45</sup> v1.0.18 to assemble reads into  
300 contigs *de novo*. Velvet Optimiser v2.1.7 was used to choose the kmer lengths on a per  
301 sequence basis. The median kmer length was 123bp (IQR 119-123). To obtain multilocus



302 sequence types we used BLAST<sup>46</sup> to find the relevant loci, and looked up the nucleotide  
303 sequences in the online database at <http://saureus.mlst.net/>. Strains that shared 6 of 7 MLST  
304 loci were considered to be in the same Clonal Complex.<sup>47</sup> Antibiotic sensitivity was predicted  
305 by interrogating the assemblies for a panel of resistance determinants as previously described.<sup>48</sup>

306 We used Stampy<sup>49</sup> v1.0.22 to map reads against reference genomes (USA300\_FPR3757,  
307 Genbank accession number CP000255.1).<sup>50</sup> Repetitive regions, defined by BLAST<sup>46</sup>  
308 comparison of the reference genome against itself, were masked prior to variant calling. Bases  
309 were called at each position using previously described quality filters.<sup>51-53</sup>

310 After filtering ONT reads with filtlong v.0.2.0 (with settings filtlong -- min\_length 1000 --  
311 keep\_percent 90 --target\_bases 500000000 --trim --split 500), hybrid assembly of short  
312 (Illumina) and long (ONT) reads were made, using Unicycler v0.4.5<sup>54</sup> (default settings). The  
313 workflow for these assemblies is available at  
314 <https://gitlab.com/ModernisingMedicalMicrobiology/MOHAWK>)

315 **Reconstructing the phylogenetic tree.** We constructed a maximum likelihood phylogeny of  
316 mapped genomes for visualization using RAxML<sup>55</sup> assuming a general time reversible (GTR)  
317 model. To overcome a limitation in the presence of divergent sequences whereby RAxML fixes  
318 a minimum branch length that may be longer than a single substitution event, we fine-tuned the  
319 estimates of branch lengths using ClonalFrameML.<sup>56</sup>

320 **Kmer counting.** We used a kmer-based approach to capture non-SNP variation.<sup>1</sup> Using the *de*  
321 *novo* assembled genome, all unique 31 base haplotypes were counted using dsk<sup>57</sup>. If a kmer  
322 was found in the assembly it was counted present for that genome, otherwise it was treated as  
323 absent. This produced a set of 10,744,013 variably present kmers, with the presence or absence  
324 of each determined per isolate. We identified a median of 2,801,000 kmers per isolate,  
325 including variably present kmers and kmers common to all genomes (IQR 2,778,000-  
326 2,837,000).

327 **Calculating heritability.** We used the Genome-wide Efficient Mixed Model Association tool  
328 (GEMMA<sup>31</sup>) to fit a univariate linear mixed model for association between a single phenotype  
329 (pyomyositis vs asymptomatic nasal carriage). We calculated the relatedness matrix from  
330 kmers, and used GEMMA to estimate the proportion of variance in phenotypes explained by  
331 genotypic diversity in the sample set (i.e. estimated heritability). Heritability estimates with and  
332 without a covariate (e.g. the presence of high risk kmers) are compared by testing for difference  
333 in proportions. We use the point estimate for heritability as the denominator to calculate the  
334 relative decrease proportion.

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

346 To address the possibility of differing effect sizes between the two control cohorts, we have  
347 repeated the analysis after splitting the study into two groups – early (2008 and earlier, n = 276,

348 cases n =54, controls n = 222) and late (2009 and later, n=242, cases n=47, controls n=195).  
349 We then examined the maximum likelihood estimates produced by the LMM for kmers  
350 mapping to the PVL coding sequence in each region. The 95% CI of the estimate from each sub  
351 study were overlapping (See Figure 2-figure supplement 4).

352 **Testing for lineage effects.** We tested for associations between lineage and phenotype using an  
353 R package *bugwas* (available at <https://github.com/sgearle/bugwas>), which implements a  
354 published method for lineage testing in bacterial GWAS.<sup>3</sup> We tested lineages using principal  
355 components. These were computed based on biallelic SNPs using the R function *prcomp*. To  
356 test the null hypothesis of no background effect of each principal component, we used a Wald  
357 test, which we compared against a  $\chi^2$  distribution with one degree of freedom to obtain a *p*  
358 value.

359 **Kmer mapping.** We used Bowtie<sup>58</sup> to align all 31bp kmers from short-read sequencing were to  
360 a draft reference (the *de novo* assembly of a CC-121 pyomyositis strain PYO2134). Areas of  
361 homology between the draft reference and well-annotated reference strains were identified by  
362 aligning sequences with Mauve.<sup>59</sup> For all 31bp kmers with significant association with case-  
363 controls status, the likely origin of the kmer was determined by nucleotide sequence BLAST<sup>46</sup>  
364 of the kmers against a database of all *S. aureus* sequences in Genbank.

365 **Joint short-read and long-read analysis.** 31bp kmers were counted for the 37 hybrid short-read  
366 and long-read assemblies using *dsk*<sup>57</sup>. The presence or absence of all Illumina (short-read)  
367 kmers that mapped to the two PVL toxin-coding sequences and the upstream intergenic region  
368 plus the surrounding 1kb were reassessed. For the 37 samples with hybrid assemblies, the  
369 presence/absence of these kmers was determined from the kmers counted from the hybrid  
370 assemblies. For all other samples, presence/absence was determined from the kmers counted  
371 from the short-read only assemblies. The new presence/absence patterns were tested for  
372 association with the phenotype controlling for population structure and genetic background  
373 using GEMMA<sup>31</sup>, using the same relatedness matrix as the original short-read analysis.

374 **Predicting presence of PVL genes and bacteriophages.** We used BLAST to check for the  
375 relative coverage of the PVL CDS (From reference genome USA300\_FPR3757 (CP000255.1)  
376 positions 1546170-1548350), as well as the entire sequence of 6 known PVL positive phages  
377 ( $\phi$ 2958(NC\_011344.1),  $\phi$ PVL (NC\_002321.1),  $\phi$ PVL108 (NC\_008689.1),  $\phi$ SLT  
378 (NC\_002661.2),  $\phi$ Sa2MW (NC\_003923.1),  $\phi$ Sa2USA (CP000255.1))<sup>18</sup>, as well as the  
379 hypothesised prophage region from PYO2134 (1571177- 1616957), which we have  
380 called  $\phi$ CC121. For PVL genes, we determined relative coverage of the query sequence; over  
381 98% coverage was used as threshold for gene presence.

382 **Multiple testing correction.** Multiple testing was accounted for by applying a Bonferroni  
383 correction;<sup>60</sup> the individual locus effect of a variant (kmer or PC) was considered significant if  
384 its *P* value was smaller than  $\alpha/n_p$ , where we took  $\alpha = 0.05$  to be the genome-wide false-positive  
385 rate and  $n_p$  to be the number of kmers or PCs with unique phylogenetic patterns, that is, unique  
386 partitions of individuals according to allele membership. We identified 236627 unique kmer  
387 patterns and 518 PCs, giving thresholds of  $2.1 \times 10^{-7}$  and  $9.7 \times 10^{-5}$  respectively.

388 **Data availability.** Sequence data has been submitted to Short Read Archive (Bioproject ID  
389 PRJNA418899). Clinical origins of sequenced strains are listed in supplementary information  
390 (Supplementary File 4).

391  
392

393 **Acknowledgments:** The authors would like to thank study participants. This study was funded  
394 by the Wellcome Trust (MORU Grants 089275/H/09/Z and 089275/Z/09/Z), and a University  
395 of Oxford Medical Research Fund awarded to C.E.M. (MRF/MT2015/2180). D.J.W. is a Sir  
396 Henry Dale Fellow, jointly funded by the Wellcome Trust and the Royal Society (Grant  
397 101237/Z/13/Z). B.C.Y. is a Research Training Fellow funded by the Wellcome Trust (Grant  
398 101611/Z/13/Z). D.H.W. was funded by the National Institute for Health Research (NIHR)  
399 Oxford Biomedical Research Centre (BRC) and the European Union's Seventh Framework  
400 Programme under the grant agreement number 601783 (BELLEROPHON project). N.S. is  
401 funded by a Public Health England (PHE)/University of Oxford Clinical Lectureship. K.E. was  
402 funded by an academic clinical fellowship which was provided by the UK NIHR through the  
403 University of Oxford. This research was supported by Core funding to the Wellcome Centre for  
404 Human Genetics provided by the Wellcome (090532/Z/09/Z). The views expressed are those of  
405 the author(s) and not necessarily those of the NHS, PHE, the NIHR or the Department of  
406 Health.

407  
408 **Competing interests:** None to declare

409  
410

411 **References:**

- 412 1. Sheppard SK, Didelot X, Meric G, Torralbo A, Jolley KA, Kelly DJ, Bentley SD, Maiden MC,  
 413 Parkhill J, Falush D. Genome-wide association study identifies vitamin B5 biosynthesis as a host  
 414 specificity factor in *Campylobacter*. *Proc Natl Acad Sci U S A*. **110(29)**, 11923-7 (2013) doi:  
 415 10.1073/pnas.1305559110.
- 416 2. Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, Mather AE, Hanage WP,  
 417 Goldblatt D, Nosten FH, Turner C, Turner P, Bentley SD, Parkhill J. Comprehensive identification  
 418 of single nucleotide polymorphisms associated with beta-lactam resistance within pneumococcal  
 419 mosaic genes. *PLoS Genet*. **10(8)**, e1004547 (2014) doi: 10.1371/journal.pgen.1004547.
- 420 3. Earle SG, Wu CH, Charlesworth J, Stoesser N, Gordon NC, Walker TM, Spencer CCA, Iqbal Z,  
 421 Clifton DA, Hopkins KL, Woodford N, Smith EG, Ismail N, Llewelyn MJ, Peto TE, Crook DW,  
 422 McVean G, Walker AS, Wilson DJ. Identifying lineage effects when controlling for population  
 423 structure improves power in bacterial association studies. *Nat Microbiol*. **1**, 16041 (2016) doi:  
 424 10.1038/nmicrobiol.2016.41.
- 425 4. Lees JA, Vehkala M, Välimäki N, Harris SR, Chewapreecha C, Croucher NJ, Marttinen P, Davies  
 426 MR, Steer AC, Tong SY, Honkela A, Parkhill J, Bentley SD, Corander J. Sequence element  
 427 enrichment analysis to determine the genetic basis of bacterial phenotypes. *Nat Commun*. **7**, 12797  
 428 (2016) doi: 10.1038/ncomms12797.
- 429 5. Falush D. Bacterial genomics: Microbial GWAS coming of age. *Nat Microbiol*. **1**, 16059 (2016) doi:  
 430 10.1038/nmicrobiol.2016.59.
- 431 6. Power R, Parkhill J, de Oliveira T. Microbial Genome-Wide Association Studies: Lessons from  
 432 Human GWAS, *Nat Rev Genet*. **18**, 41-50 (2017) doi: 10.1038/nrg.2016.132.
- 433 7. Chauhan S, Jain S, Varma S, Chauhan SS. Tropical pyomyositis (myositis tropicans): current  
 434 perspective. *Postgrad Med J*. **80(943)**, 267-70. (2004).
- 435 8. Verma S. Pyomyositis in Children. *Curr Infect Dis Rep*. **18(4)**, 12 (2016) doi: 10.1007/s11908-016-  
 436 0520-2.
- 437 9. Bickels J, Ben-Sira L, Kessler A, Wientroub S. Primary pyomyositis. *J Bone Joint Surg Am*. **84-**  
 438 **A(12)**, 2277-86 (2002).
- 439 10. Moriarty P, Leung C, Walsh M, Nourse C. Increasing pyomyositis presentations among children in  
 440 Queensland, Australia. *Pediatr Infect Dis J*. **34(1)**, 1-4 (2015) doi: 10.1097/INF.0000000000000470.
- 441 11. Pannaraj PS, Hulten KG, Gonzalez BE, Mason EO Jr, Kaplan SL. Infective pyomyositis and  
 442 myositis in children in the era of community-acquired, methicillin-resistant *Staphylococcus aureus*  
 443 infection. *Clin Infect Dis*. **43(8)**, 953-60 (2006).
- 444 12. Borges AH, Faragher B, Laloo DG. Pyomyositis in the upper Negro river basin, Brazilian  
 445 Amazonia. *Trans R Soc Trop Med Hyg*. **106(9)**, 532-7 (2012) doi: 10.1016/j.trstmh.2012.06.008.
- 446 13. Shallcross LJ, Fragaszy E, Johnson AM, Hayward AC. The role of the Pantone-Valentine leucocidin  
 447 toxin in staphylococcal disease: a systematic review and meta-analysis. *Lancet Infect Dis*. **13**, 43-54  
 448 (2013) doi: 10.1016/S1473-3099(12)70238-4.
- 449 14. Vandenesch F, Naimi T, Enright MC, Lina G, Nimmo GR, Heffernan H, Liassine N, Bes M,  
 450 Greenland T, Reverdy ME, Etienne J. Community-acquired methicillin-resistant *Staphylococcus*  
 451 *aureus* carrying Pantone-Valentine leukocidin genes: worldwide emergence. *Emerg Infect Dis*. **9(8)**,  
 452 978-84 (2003).
- 453 15. Labandeira-Rey M, Couzon F, Boisset S, Brown EL, Bes M, Benito Y, Barbu EM, Vazquez V,  
 454 Höök M, Etienne J, Vandenesch F, Bowden MG. *Staphylococcus aureus* Pantone-Valentine  
 455 leukocidin causes necrotizing pneumonia. *Science*. **315(5815)**, 1130-3 (2007).
- 456 16. Villaruz AE, Bubeck Wardenburg J, Khan BA, Whitney AR, Sturdevant DE, Gardner DJ, DeLeo  
 457 FR, Otto M. A point mutation in the agr locus rather than expression of the Pantone-Valentine

- 458 leukocidin caused previously reported phenotypes in *Staphylococcus aureus* pneumonia and gene  
459 regulation. *J Infect Dis.* **200(5)**, 724-34 (2009) doi: 10.1086/604728.
- 460 17. Löffler B, Hussain M, Grundmeier M, Brück M, Holzinger D, Varga G, Roth J, Kahl BC, Proctor  
461 RA, Peters G. *Staphylococcus aureus* Panton-Valentine leukocidin is a very potent cytotoxic factor  
462 for human neutrophils. *PLoS Pathog.* **6(1)**, e1000715 (2010) doi: 10.1371/journal.ppat.1000715.
- 463 18. Boakes E, Kearns AM, Ganner M, Perry C, Hill RL, Ellington MJ. Distinct bacteriophages encoding  
464 Panton-Valentine leukocidin (PVL) among international methicillin-resistant *Staphylococcus aureus*  
465 clones harboring PVL. *J Clin Microbiol.* **49(2)**, 684-92 (2011) doi: 10.1128/JCM.01917-10.
- 466 19. McCarthy AJ, Witney AA, Lindsay JA. *Staphylococcus aureus* temperate bacteriophage: carriage  
467 and horizontal gene transfer is lineage associated. *Front Cell Infect Microbiol.* **2**, 6 (2012) doi:  
468 10.3389/fcimb.2012.00006.
- 469 20. Oliveira D, Borges A, Simões M. *Staphylococcus aureus* Toxins and Their Molecular Activity in  
470 Infectious Diseases. *Toxins (Basel).* **10(6)**, pii: E252. (2018) doi: 10.3390/toxins10060252.
- 471 21. Niemann S, Bertling A, Brodde MF, et al Panton-Valentine Leukocidin associated with *S. aureus*  
472 osteomyelitis activates platelets via neutrophil secretion products. *Sci Rep.* **8(1)**, 2185. (2018) doi:  
473 10.1038/s41598-018-20582-z.
- 474 22. Spaan AN, Henry T, van Rooijen WJM, et al. The staphylococcal toxin Panton-Valentine  
475 Leukocidin targets human C5a receptors. *Cell Host Microbe.* **13(5)**, 584-594. (2013) doi:  
476 10.1016/j.chom.2013.04.006
- 477 23. Sina H, Ahoyo TA, Moussaoui W, Keller D, Bankolé HS, Barogui Y, Stienstra Y, Kotchoni SO,  
478 Prévost G, Baba-Moussa L. Variability of antibiotic susceptibility and toxin production of  
479 *Staphylococcus aureus* strains isolated from skin, soft tissue, and bone related infections. *BMC*  
480 *Microbiol.* **13**, 188 (2013) doi:10.1186/1471-2180-13-188.
- 481 24. García C, Hallin M, Deplano A, Denis O, Sihuinchu M, de Groot R, Gotuzzo E, Jacobs J.  
482 *Staphylococcus aureus* causing tropical pyomyositis, Amazon Basin, Peru. *Emerg Infect Dis.* **19(1)**,  
483 123-5 (2013) doi: 10.3201/eid1901.120819.
- 484 25. Otto M. A MRSA-terious enemy among us: end of the PVL controversy? *Nat Med.* **17(2)**, 169-70  
485 (2011) doi: 10.1038/nm0211-169.
- 486 26. Day NPJ. Panton-Valentine leucocidin and staphylococcal disease *Lancet Infect Dis.* **13**: 5-6 (2013)  
487 doi: 10.1016/S1473-3099(12)70265-7.
- 488 27. Lees JA, Croucher NJ, Goldblatt D, Nosten F, Parkhill J, Turner C, Turner P, Bentley SD. Genome-  
489 wide identification of lineage and locus specific variation associated with pneumococcal carriage  
490 duration. *Elife.* **6**; pii: e26255 (2017) doi: 10.7554/eLife.2625.
- 491 28. Méric G, Mageiros L, Pensar J, Laabei M, Yahara K, Pascoe B, Kittivan N, Tadee P, Post V,  
492 Lamble S, Bowden R, Bray JE, Morgenstern M, Jolley KA, Maiden MCJ, Feil EJ, Didelot X,  
493 Miragaia M, de Lencastre H, Moriarty TF, Rohde H, Massey R, Mack D, Corander J, Sheppard SK.  
494 Disease-associated genotypes of the commensal skin bacterium *Staphylococcus epidermidis*. *Nat*  
495 *Commun.* **9(1)**: 5034 (2018) doi: 10.1038/s41467-018-07368-7.
- 496 29. Chewapreecha C, Holden MT, Vehkala M, Välimäki N, Yang Z, Harris SR, Mather AE, Tuanyok A,  
497 De Smet B, Le Hello S, Bizet C, Mayo M, Wuthiekanun V, Limmathurotsakul D, Phetsouvanh R,  
498 Spratt BG, Corander J, Keim P, Dougan G, Dance DA, Currie BJ, Parkhill J, Peacock SJ. Global and  
499 regional dissemination and evolution of *Burkholderia pseudomallei*. *Nat Microbiol.* **2**:16263. (2017)  
500 doi: 10.1038/nmicrobiol.2016.263.
- 501 30. Recker M, Laabei M, Toleman MS, Reuter S, Saunderson RB, Blane B, Török ME, Ouadi K,  
502 Stevens E, Yokoyama M, Steventon J, Thompson L, Milne G, Bayliss S, Bacon L, Peacock SJ,  
503 Massey RC. Clonal differences in *Staphylococcus aureus* bacteraemia-associated mortality. *Nat*  
504 *Microbiol.* **2(10)**:1381-1388. (2017) doi: 10.1038/s41564-017-0001-x

- 505 31. Zhou X and Stephens M. Genome-wide efficient mixed-model analysis for association studies.  
506 *Nature Genetics*. **44**, 821–824 (2012) doi: 10.1038/ng.2310.
- 507 32. Holden MT, Feil EJ, Lindsay JA, et al. Complete genomes of two clinical *Staphylococcus aureus*  
508 strains: evidence for the rapid evolution of virulence and drug resistance. *Proc Natl Acad Sci U S A*.  
509 **101(26)**:9786-9 (2004).
- 510 33. Xia G, Wolz C. Phages of *Staphylococcus aureus* and their impact on host evolution. *Infect Genet*  
511 *Evol*. **21**:593-601. (2014) doi: 10.1016/j.meegid.2013.04.022
- 512 34. M. Stegger, T. Wirth, P.S. Andersen, R.L.Skov, A. De Grassi, P.M. Simões, et al. Origin and  
513 evolution of European community-acquired methicillin-resistant *Staphylococcus aureus*  
514 *MBio*, **5** (2014) e01044-14. doi: 10.1128/mBio.01044-14.
- 515 35. Diep BA, Palazzolo-Ballance AM, Tattévin P, Basuino L, Braughton KR, Whitney AR, Chen L,  
516 Kreiswirth BN, Otto M, DeLeo FR, Chambers HF. Contribution of Panton-Valentine leukocidin in  
517 community-associated methicillin-resistant *Staphylococcus aureus* pathogenesis. *PLoS One*. **3(9)**,  
518 e3198 (2008) doi: 10.1371/journal.pone.0003198.
- 519 36. Bocchini CE, Hultén KG, Mason EO Jr, Gonzalez BE, Hammerman WA, Kaplan SL. Panton-  
520 Valentine leukocidin genes are associated with enhanced inflammatory response and local disease in  
521 acute hematogenous *Staphylococcus aureus* osteomyelitis in children. *Pediatrics*. **117(2)**, 433-40  
522 (2006).
- 523 37. Kurt K, Rasigade JP, Laurent F, Goering RV, Žemličková H, Machova I, Struelens MJ, Zautner AE,  
524 Holtfreter S, Bröker B, Ritchie S, Reaksmey S, Limmathurotsakul D, Peacock SJ, Cuny C, Layer F,  
525 Witte W, Nübel U. Subpopulations of *Staphylococcus aureus* Clonal Complex 121 Are Associated  
526 with Distinct Clinical Entities. *PLoS ONE*. **8(3)**, e58155 (2013) doi: 10.1371/journal.pone.0058155.
- 527 38. Nickerson EK, Wuthiekanun V, Kumar V, Amornchai P, Wongdeethai N, Chheng K, Chantratita N,  
528 Putschhat H, Thaipadungpanit J, Day NP, Peacock SJ. Emergence of community-associated  
529 methicillin-resistant *Staphylococcus aureus* carriage in children in Cambodia. *Am J Trop Med Hyg*.  
530 **84(2)**, 313-7 (2011) doi: 10.4269/ajtmh.2011.10-0300.
- 531 39. Miyake H. Beitragezurkenntnis der sogenannten myositis infectiosa. *Mitt Grenageb Med Chir*.  
532 **13**:155–98. (1904).
- 533 40. Tseng CW, Kyme P, Low J, Rocha MA, Alsabeh R, Miller LG, Otto M, Arditi M, Diep BA, Nizet  
534 V, Doherty TM, Beenhouwer DO, Liu GY. *Staphylococcus aureus* Panton-Valentine leukocidin  
535 contributes to inflammation and muscle tissue injury. *PLoS One*.**4(7)**:e6387. (2009) doi:  
536 10.1371/journal.pone.0006387.
- 537 41. Lehman D, Tseng CW, Eells S, Miller LG, Fan X, Beenhouwer DO, Liu GY. *Staphylococcus aureus*  
538 Panton-Valentine leukocidin targets muscle tissues in a child with myositis and necrotizing fasciitis.  
539 *Clin Infect Dis*.**50(1)**:69-72. (2010) doi: 10.1086/649217.
- 540 42. Landrum ML, Lalani T, Niknian M, Maguire JD, Hospenthal DR, Fattom A, Taylor K, Fraser J,  
541 Wilkins K, Ellis MW, Kessler PD, Fahim RE, Tribble DR. Safety and immunogenicity of a  
542 recombinant *Staphylococcus aureus*  $\alpha$ -toxoid and a recombinant Panton-Valentine leukocidin  
543 subunit, in healthy adults. *Hum Vaccin Immunother*.**13(4)**, 791-801 (2017) doi:  
544 10.1080/21645515.2016.1248326.
- 545 43. Saeed K, Gould I, Esposito S, Ahmad-Saeed N, Ahmed SS, Alp E, Bal AM, Bassetti M, Bonnet E,  
546 Chan M, Coombs G, Dancer SJ, David MZ, De Simone G, Dryden M, Guardabassi L, Hanitsch LG,  
547 Hijazi K, Krüger R, Lee A, Leistner R, Pagliano P, Righi E, Schneider-Burrus S, Skov RL, Tattévin  
548 P, Van Wamel W, Vos MC, Voss A; International Society of Chemotherapy. Panton-Valentine  
549 leukocidin-positive *Staphylococcus aureus*: a position statement from the International Society of  
550 Chemotherapy. *Int J Antimicrob Agents*. 2018 Jan;**51(1)**:16-25. doi:  
551 10.1016/j.ijantimicag.2017.11.002.
- 552 44. CLSI. *Performance Standards for Antimicrobial Susceptibility Testing*. 24th ed. CLSI supplement  
553 M100. Wayne, PA: Clinical and Laboratory Standards Institute; 2014.

- 554 45. Zerbino DR and Birney E. Velvet: Algorithms for de novo short read assembly using de bruijn  
555 graphs *Genome Res.* **18(5)**, 821-9 (2008) doi: 10.1101/gr.074492.107.
- 556 46. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool *J Mol*  
557 *Biol.* **215(3)**, 403-10 (1990).
- 558 47. Feil EJ, Cooper JE, Grundmann H, Robinson DA, Enright MC, Berendt T, Peacock SJ, Smith JM,  
559 Murphy M, Spratt BG, Moore CE, Day NP. How clonal is *Staphylococcus aureus*? *J Bacteriol.*  
560 **185(11)**:3307-16 (2003).
- 561 48. Gordon NC, Price JR, Cole K, Everitt R, Morgan M, Finney J, Kearns AM, Pichon B, Young B,  
562 Wilson DJ, Llewelyn MJ, Paul J, Peto TE, Crook DW, Walker AS, Golubchik T. Prediction of  
563 *Staphylococcus aureus* antimicrobial resistance by whole-genome sequencing *J Clin Microbiol.*  
564 **52(4)**,1182-91 (2014) doi: 10.1128/JCM.03117-13.
- 565 49. Lunter G and Goodson M. Stampy: A statistical algorithm for sensitive and fast mapping of illumina  
566 sequence reads *Genome Res.* **21(6)**, 936-9 (2011) doi: 10.1101/gr.111120.110.
- 567 50. Diep BA, Gill SR, Chang RF, Phan TH, Chen JH, Davidson MG, Lin F, Lin J, Carleton HA,  
568 Mongodin EF, Sensabaugh GF, Perdreau-Remington F. Complete genome sequence of USA300, an  
569 epidemic clone of community-acquired methicillin-resistant *Staphylococcus aureus*. *Lancet.*  
570 **367(9512)**, 731-9 (2006).
- 571 51. Didelot X, Eyre DW, Cule M, Ip CL, Ansari MA, Griffiths D, Vaughan A, O'Connor L, Golubchik  
572 T, Batty EM, Piazza P, Wilson DJ, Bowden R, Donnelly PJ, Dingle KE, Wilcox M, Walker AS,  
573 Crook DW, Peto TE, Harding RM. Microevolutionary analysis of *Clostridium difficile* genomes to  
574 investigate transmission. *Genome Biol.* **13(12)**,R118 (2012) doi: 10.1186/gb-2012-13-12-r118.
- 575 52. Young BC, Golubchik T, Batty EM, Fung R, Larner-Svensson H, Votintseva AA, Miller RR,  
576 Godwin H, Knox K, Everitt RG, Iqbal Z, Rimmer AJ, Cule M, Ip CL, Didelot X, Harding RM,  
577 Donnelly P, Peto TE, Crook DW, Bowden R, Wilson DJ. Evolutionary dynamics of *Staphylococcus*  
578 *aureus* during progression from carriage to disease. *Proc Natl Acad Sci U S A.* **109(12)**, 4550 (2012)  
579 doi: 10.1073/pnas.1113219109.
- 580 53. Golubchik T, Batty EM, Miller RR, Farr H, Young BC, Larner-Svensson H, Fung R, Godwin H,  
581 Knox K, Votintseva A, Everitt RG, Street T, Cule M, Ip CL, Didelot X, Peto TE, Harding RM,  
582 Wilson DJ, Crook DW, Bowden R. Within- host evolution of *Staphylococcus aureus* during  
583 asymptomatic carriage. *PloS One.* **8(5)**: e61319 (2013) doi: 10.1371/journal.pone.0061319.
- 584 54. Wick RR, Judd LM, Gorrie CL, Holt KE. Unicycler: Resolving bacterial genome assemblies from  
585 short and long sequencing reads. *PLoS Comput Biol.* **13(6)**:e1005595. (2017) doi:  
586 10.1371/journal.pcbi.1005595.
- 587 55. Stamatakis A. RAxML version 8: A tool for phylogenetic analysis and post-analysis of large  
588 phylogenies *Bioinformatics.* **30(9)**,1312-3 (2014) doi: 10.1093/bioinformatics/btu033.
- 589 56. Didelot X and Wilson DJ. ClonalFrameML: Efficient inference of recombination in whole bacterial  
590 genomes *PLoS Comput Biol.* **11(2)**, e1004041 (2015) doi: 10.1371/journal.pcbi.1004041.
- 591 57. Rizk G, Lavenier D, and Chikhi R. DSK: k-mer counting with very low memory usage.  
592 *Bioinformatics.* **29**, 652–653 (2013) doi: 10.1093/bioinformatics/btt020.
- 593 58. Langmead B, Salzberg SL. Fast gapped-read alignment with bowtie 2. *Nat Methods.* **4**; 9(4), 357-9  
594 (2012) doi: 10.1038/nmeth.1923.
- 595 59. Darling AC, Mau B, Blattner FR, Perna NT. Mauve: Multiple alignment of conserved genomic  
596 sequence with rearrangements. *Genome Res.* **14(7)**, 1394-403 (2004) doi: 10.1101/gr.2289704
- 597 60. Dunn OJ. Estimation of the medians for dependent variables. *Ann. Math. Stat.* **30**, 192–197 (1959).

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

608 **Figure 1-figure supplement 1** Sampling frequencies of the major strains were stable over  
609 time. The year of sampling (2007-2008, blue shaded lines) and MRSA status (orange lines)  
610 are illustrated around the phylogeny of the bacteria sampled from pyomyositis cases (gold  
611 lines) and asymptomatic carriage controls (green lines). The three PCs most significantly  
612 associated with case/control status are also shown (PCs 1, 2 and 20 by red, blue and yellow  
613 branches respectively). Branch lengths are square root transformed to aid visualization.

614 **Figure 2.** Kmers associated with pyomyositis. **(A)** All kmers (n = 10,744,013) were mapped  
615 to the genome assembly of one CC121 pyomyositis bacterium (PYO2134). Each point  
616 represents a kmer, plotted by the mapped location and the significance of the association with  
617 disease ( $-\log_{10} p$  value). Kmers are coloured by the odds ratio (OR) of kmer presence for  
618 disease risk. A Bonferroni-adjusted threshold for significance is plotted in grey **(B)** The  
619 region between 1.57-1.62 MB in greater detail. Grey arrows depict coding sequences,  
620 determined by homology to USA300 FPR3757. **(C)** Associations for kmers mapping to  
621 region 1,571 – 1,574kB is plotted. Kmer presence determined from hybrid assembly using  
622 short and long-read data for assembly. Grey arrows depict coding sequences, determined by  
623 homology to USA300 FPR3757.

624  
625 **Figure 2-figure supplement 1:** Alignments of reference genome PYO2134 assembly (R)  
626 with 37 *de novo* assemblies of Illumina short-read sequencing (C) which feature either  
627 ambiguities (Ns) or contig boundaries in the region 389bp upstream of PVL coding sequence.  
628 Contig boundaries, when overlapping, are marked with a red diamond. Ns in the assembly are  
629 marked in dark grey. Polymorphisms are colour-coded by base.

630  
631 **Figure 2-figure supplement 2:** Alignments of reference genome PYO2134 assembly (R)  
632 with 37 *de novo* assemblies of Illumina short-read sequencing (C) which feature either  
633 ambiguities (Ns) or contig boundaries in the region 389bp upstream of PVL coding sequence.  
634 Contig boundaries, when overlapping, are marked with a red diamond. Ns in the assembly are  
635 marked in dark grey. Polymorphisms are colour-coded by base.

636 **Figure 2-figure supplement 3:** Presence of PVL and potential PVL carrying phages across  
637 the population. Multi-locus sequence type (ST) or clonal complex (CC) groups are shown  
638 (outer grey ring). Case/control status is marked in the next outermost ring: pyomyositis (gold,  
639 n = 101) or nasal carriage (green, n = 417). Branch lengths are square root transformed to aid  
640 visualization. The presence of kmers mapping to the PVL coding sequence (dark red) or of  
641 >98% coverage of PVL coding sequence in *de novo* assembly found by BLAST (red) are  
642 marked in the next rings. Phages with >95% coverage on BLAST of the *de novo* assembly  
643 are marked by points ( $\phi$ PVL in purple,  $\phi$ SaUSA2 in dark blue,  $\phi$ SLT in light blue and



644  $\phi$ CC121 in orange), and bars are used to mark the phage with the greatest sequence  
645 homology.

646

647 **Figure 2-figure supplement 4:** Significant association between case/control status and the  
648 presence of kmers mapping to PVL was consistent across early and late subsets. Results  
649 from repeat kmer GWAS using early and late sub-groups of cases and controls, as well as  
650 results from original GWAS of the entire group. We report the presence of kmers mapping  
651 across the PVL locus in each subset. Forest plot showing the effect sizes (box) with 95% CI  
652 (whiskers) determined by LMM for kmers mapping to PVL kmer in each subset, and in the  
653 complete study (diamond centred on effect size extending across 95% CI). The statistical  
654 significance of association with case/control status after controlling for population structure is  
655 shown for each subset and the complete study.

656

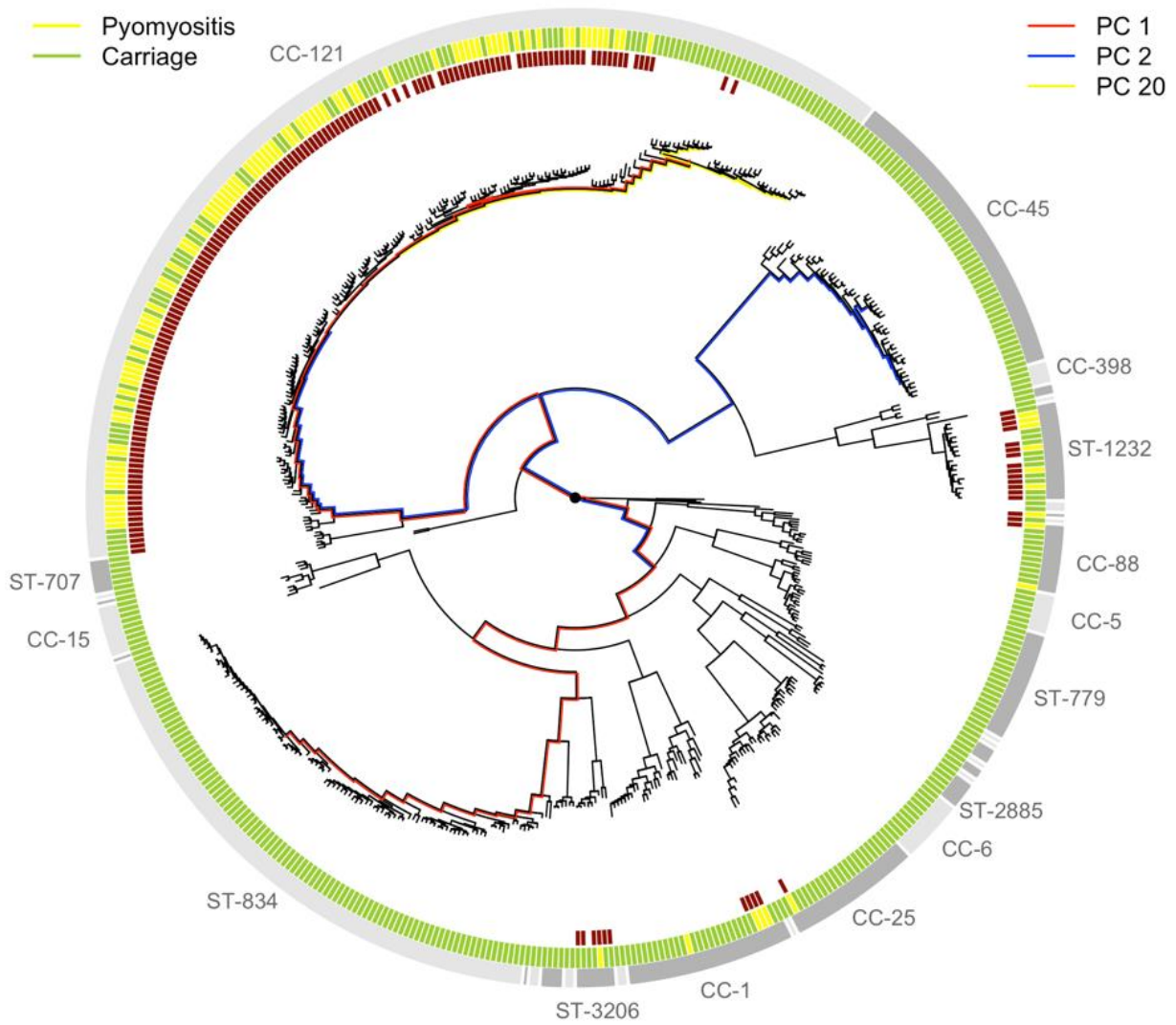
657 **Supplementary File 1:** Isolates included in this study.

658 **Supplementary File 2:** All significant kmers from short read sequencing assembly, evidence  
659 of association, frequency, location on mapping to the study reference PYO2134, best match  
660 on BLAST (blastn) to all *S. aureus* coding sequences in Genbank, and best match on BLAST  
661 (blastn) to the NCBI database.

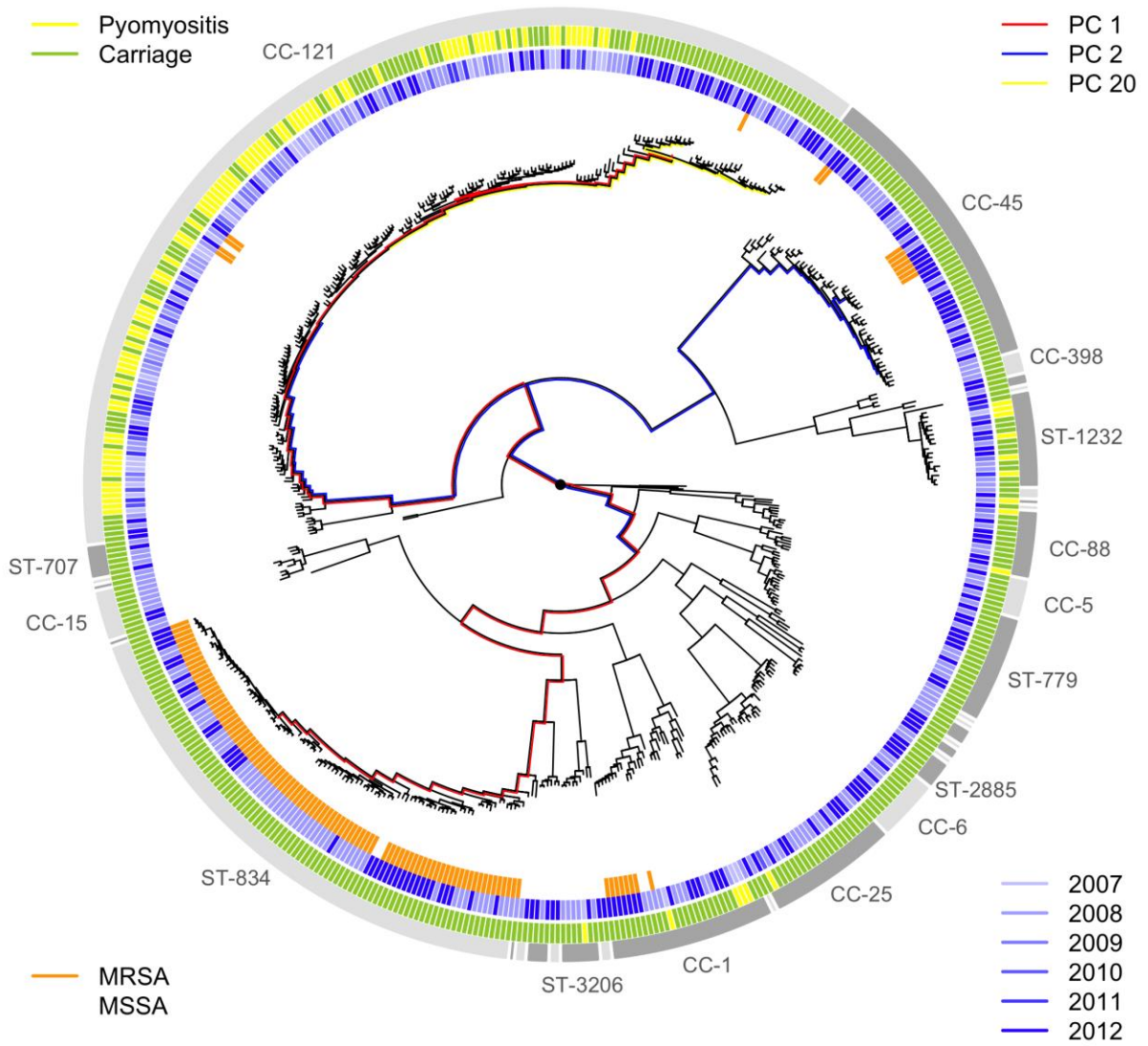
662 **Supplementary File 3:** Presence of high risk kmers and relative coverage of PVL coding  
663 sequence and common PVL positive phages found by BLAST (blastn) of short read  
664 assemblies.

665 **Supplementary File 4:** List of all isolates, and site of isolation (carriage or invasive disease)  
666 and year of isolation. These isolate names match those used in sequence data deposition on  
667 SRA.

668

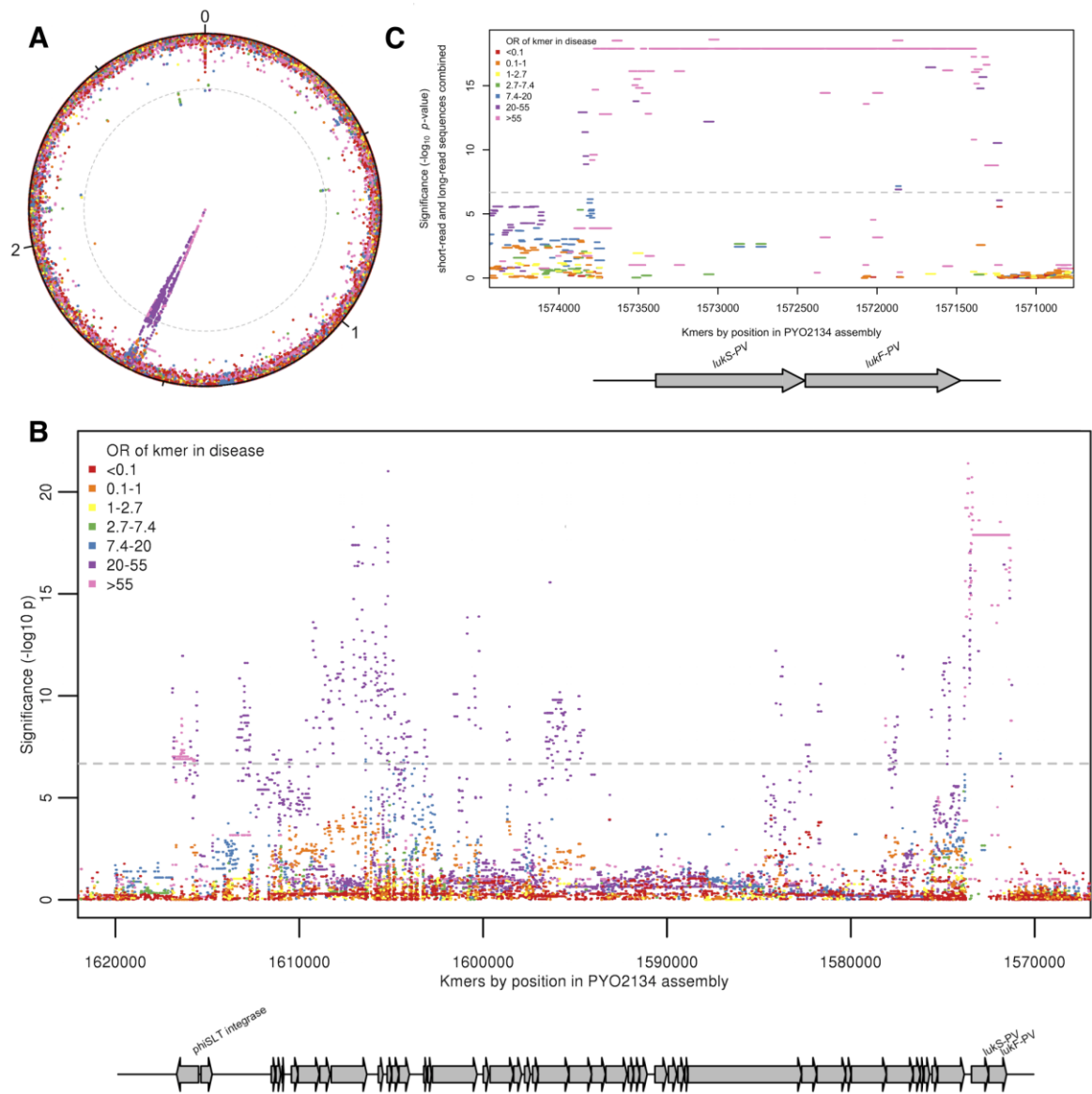


1  
2 **Figure 1**



1

2 **Figure 1-figure supplement 1**



1

2 **Figure 2**

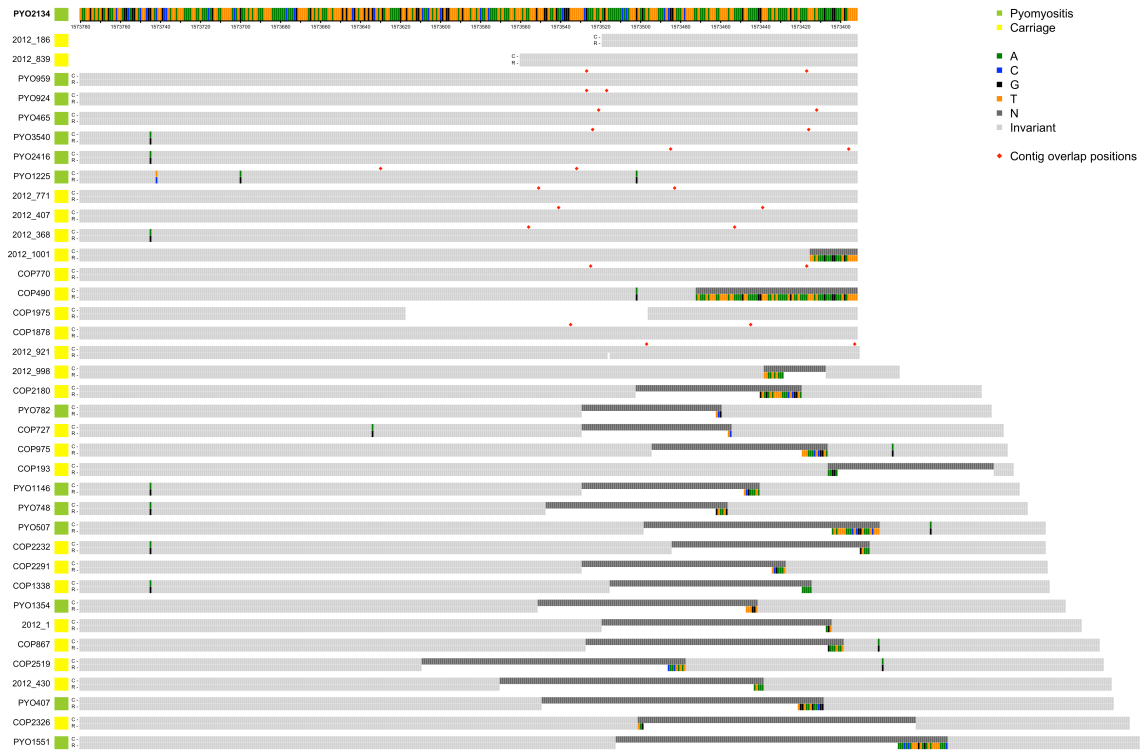


Figure 2-figure supplement 1

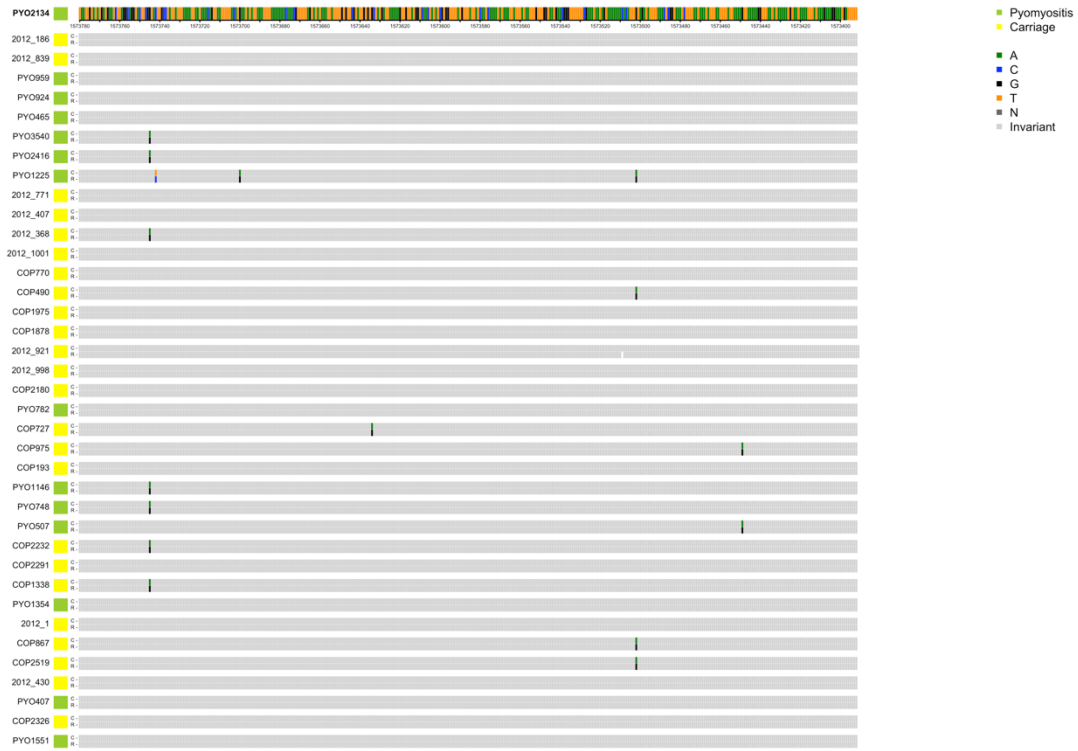
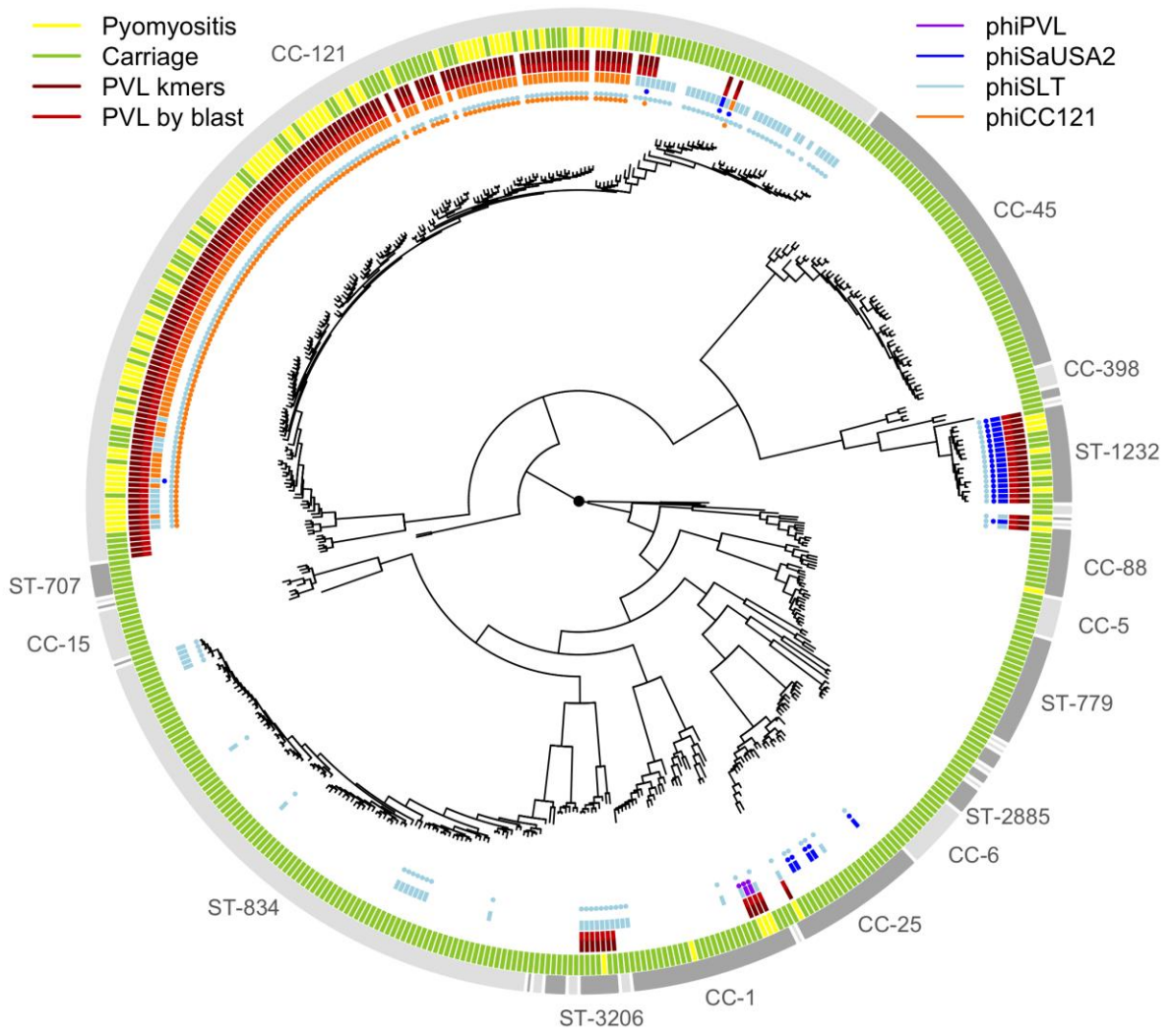


Figure 2-figure supplement 2



**Figure 2-figure supplement 3**

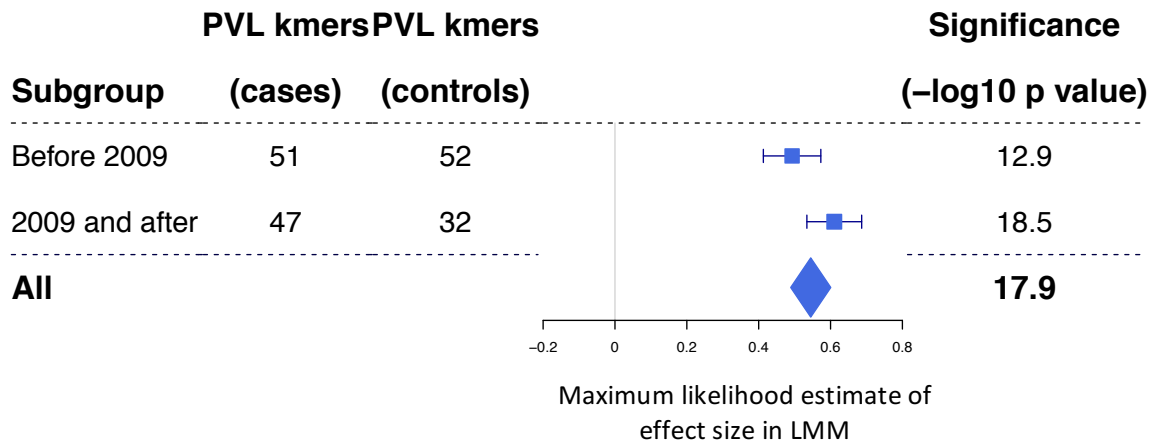


Figure 2-figure supplement 4