MICROBIOLOGY SOCIETY

1  **MinSNPs: an R package for derivation of resolution-**

2  **optimised SNP sets from microbial genomic data**

3  Kian Soon Hoon[1], Deborah C Holt[2, 1], Sarah Auburn[1,3,4], Peter Shaw[5], Philip M. Giffard[1,2]

4  **1.1  Affiliation**

5  [1]Division of Global and Tropical Health, Menzies School of Health Research, Charles Darwin

6  University, Darwin, Northern Territory, Australia.

7  [2]College of Health and Human Sciences, Charles Darwin University, Darwin, Northern Territory,

8  Australia.

9  [3]Mahidol-Oxford Tropical Medicine Research Unit, Mahidol University, Bangkok, Thailand

10  [4]Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine, University of

11  Oxford, Oxford, United Kingdom

12  [5] Oujiang Laboratory, Wenzhou, Zhejiang, China

13  Corresponding author: Philip M Giffard

14  **1.2  Keyword**

15  SNP, Bacteria, Plasmodium, Staphylococcus, resolution optimised, genotyping, surveillance, SNP

16  matrix, SNP set derivation, genetic epidemiology.

17  **1.3  Repositories:**

18  https://github.com/ludwigHoon/minSNPs

19  https://cran.r-project.org/package=minSNPs

20    https://figshare.com/s/73e279b25b1c6b59189c

21    https://microreact.org/project/minsnps-starrs

22    https://figshare.com/s/696f696c232404f18a36

23    https://figshare.com/s/746dd263140963185c53

24    https://figshare.com/s/8adc2b14052ccb89dbed

25    https://figshare.com/s/db47a069aab93f3c615c

26    https://figshare.com/s/aadf860f3cfd9c416e3f

## 2.  Abstract

28    Here we present the R package - MinSNPs.  This is designed to assemble resolution optimised sets of

29    single nucleotide polymorphisms (SNPs) from alignments such as genome wide orthologous SNP

30    matrices. We also demonstrate a pipeline for assembling such matrices from multiple bio-projects,

31    so as to facilitate SNP set derivation from globally representative data sets. MinSNPs can derive sets

32    of SNPs optimised for discriminating any user-defined combination of sequences from all others.

33    Alternatively, SNP sets may be optimised to discriminate all from all, i.e., to maximise diversity.

34    MinSNPs encompasses functions that facilitate rapid and flexible SNP mining, and clear and

35    comprehensive presentation of the results. The MinSNPs running time scales in a linear fashion with

36    input data volume, and the numbers of SNPs and SNPs sets specified in the output. MinSNPs was

37    tested using a previously reported orthologous SNP matrix of *Staphylococcus aureus*. and an

38    orthologous SNP matrix of 3,279 genomes with 164,335 SNPs assembled from four *S. aureus* short

39    read genomic data sets.  MinSNPs demonstrated efficacy in deriving discriminatory SNP sets for

40    potential surveillance targets and in identifying SNP sets optimised to discriminate isolates from

41    different clonal complexes (CC).  MinSNPs was also tested with a large *Plasmodium vivax*

42    orthologous SNP matrix. A set of five SNPs was derived that reliably indicated the country of origin

43    within 3 south-east Asian countries. In summary, we report the capacity to assemble comprehensive

44    SNP matrices that effectively capture microbial genomic diversity, and to rapidly and flexibly mine

45    these entities for optimised surveillance marker sets.


## 3. Impact statement

47    We present the R package "MinSNPs". This derives resolution optimised SNP sets from datasets of

48    genome sequence variation. Such SNP sets can underpin targeted genetic analysis for high

49    throughput surveillance of microbial variants of public health concern. MinSNPs supports

50    considerable flexibility in search methods. The package allows non-specialist bioinformaticians to

51    easily and quickly convert global scale data of intra-specific genomic variation into SNP sets precisely

52    and efficiently directed towards many microbial genetic analysis tasks.


## 4. Data summary

54       1.   The source code for minSNPs is available from GitHub under MIT Licence (URLs –

55            https://github.com/ludwigHoon/minSNPs and mirrored in https://cran.r-

56            project.org/package=minSNPs)

57       2.   *Staphylococcus aureus* (STARRS data set) Orthologous SNP Matrix; (URL -

58            https://doi.org/10.1371/journal.pone.0245790.s005)

59       3.   *Plasmodium vivax* data set (VCF file); (URL - https://www.malariagen.net/resource/24)

60       4.   *Staphylococcus aureus* short read sequences (fastq) from bioprojects: PRJEB40888 (or

61            STARRS)(https://www.ncbi.nlm.nih.gov/bioproject/PRJEB40888), PRJEB3174

62            (https://www.ncbi.nlm.nih.gov/bioproject/PRJEB3174), PRJEB32286

63            (https://www.ncbi.nlm.nih.gov/bioproject/PRJEB32286), and PRJNA400143

64            (https://www.ncbi.nlm.nih.gov/bioproject/PRJNA400143)

65    **The authors confirm all supporting data, code and protocols have been provided within the article**

66    **or through supplementary data files.**

## 5.  Introduction

68    The extremely large-scale accumulation of microbial whole genome sequence information provides

69    a potent resource for the design of targeted genetic analysis procedures. While whole genome

70    analysis is now widely applied directly to public health, clinical, and research microbiology, targeted

71    genetic analyses may be complementary to whole genome analysis for purposes such as high-

72    volume, low-cost surveillance, analysis of primary specimens, and/or analyses performed outside the

73    laboratory environment.  Several research groups have recently developed SNP-based genotyping

74    approaches, e.g., to investigate *Mycobacterium species* (1, 2), attribute host for *Chlamydia psittaci*

75    (3) and *Campylobacter coli (4)*, distinguish *Rickettsia typhi* from different continents (5), identify

76    *Escherichia coli* of specific serotype (6), and track the spread of drug resistance in *Plasmodium*

77    *falciparum* infections (7).

78    Here we report the R package "MinSNPs". This package is designed to derive sets of polymorphisms

79    from biological sequence alignment data on the basis of high combinatorial discriminatory power.

80    The envisioned application is the derivation of high-resolution sets of single nucleotide

81    polymorphisms (SNPs) from DNA sequence alignments or orthologous SNP matrices. MinSNPs

82    encompasses much of the functionality of the previously reported "Minimum SNPs" Java-based

83    bioinformatics application (8, 9). Minimum SNPs was used to develop a number of SNP-based

84    bacterial genotyping methods e.g., (10-14). MinSNPs is a new package, written in R, with distinct

85    code from Minimum SNPs. The reasons for re-development were improvement of flexibility, error

86    handling, and output formats.

87   Here we describe MinSNPs and demonstrate functionality using comparative genome data from

88   *Staphylococcus aureus* and *Plasmodium vivax.* We also demonstrate a pipeline to generate MinSNPs

89   input files from multiple short read data sets to facilitate the analysis of data from multiple studies.

## 6.   Theory & implementation

91   The input format is a single sequence alignment in FASTA format. All symbols are recognised so that

92   the program will derive sets of polymorphic positions from any file in a FASTA format alignment,

93   irrespective of the symbols in the sequences. However, symbols that are not G, A, T or C can

94   optionally trigger the exclusion of the relevant alignment positions from analysis. Our focus has

95   largely been on the analysis of genome-wide orthologous SNP matrices.

96   The output of MinSNPs is set(s) of polymorphic positions in the alignment. SNP sets are assembled

97   iteratively, on the basis of maximised combinatorial resolving power. SNP 1 is the single SNP with the

98   highest resolving power, SNP 2 is the SNP with the highest resolving power in combination with SNP

99   1 etc.  Where more than one SNP confers the same increase in resolving power, the SNP nearest to

100   position 1 will be added to the set.

101   There are two user-selectable algorithms for measuring resolving power.

102       1.    **% mode.** The resolving power is the percentage of sequences in the alignment that are

103             not discriminated from the user-selected sequence(s) (the group of interest). The SNP

104             sets are constrained to 100% sensitivity. The first SNP identified is the 100% sensitive

105             SNP with maximum possible specificity, while subsequent SNPs are selected on the basis

106             of the maximum possible increase in specificity in combination with the previously

107             selected SNP(s). All alignment positions that are variable within the group of interest can

108             optionally be excluded from the analysis. We suggest that, where possible, the group of

109         interest be composed of >1 sequence to avoid the identification of spurious SNPs arising

110         from sequencing errors.

111     2.     **D mode**.  The resolving power is the power to discriminate "all from all", as measured by

112         the Simpsons index of diversity (*D*).  In this context, *D* is the probability that any two

113         sequences in the alignment will be discriminated from each other by the SNP set, as

114         calculated by $D = 1 - \frac{1}{N(N-1)}\sum_{j=1}^{s} n_j(n_j - 1)$, where N is the number of sequences, s is

115         the number of classes defined by the SNPs, and $n_j$ is the number of sequences defined

116         by the class j (15).

117     MinSNPs encompasses functions that support flexibility of analyses and transparency of outputs. The

118     user specifies the size and number of the SNP sets that constitute the output. When multiple SNP

119     sets are requested, MinSNPs identifies alternative SNP sets that are all resolution optimised, with

120     the constraint that the sets must differ from each other at least in the first SNP.

121     The user can force the program to include or exclude any alignment position(s) in/from the SNP set.

122     Where positions are included, new SNPs are identified based on resolving power in combination

123     with the included positions. This facilitates rapid modification of SNP sets.

124     MinSNPs can identify alignment positions where at least one sequence has a non-standard DNA

125     symbol, and these positions are optionally excluded from analysis. Indels (dashes) default to being

126     regarded as symbols equivalent to other symbols. Alternatively, the user can specify that indels

127     trigger the exclusion of the relevant alignment positions from the analysis. There is also an optional

128     function to exclude positions with SNPs with >2 alleles.

129     MinSNPs provides a cumulative increase in resolving power as the sets are built, and the tabulated

130     information indexing the sequences in the alignment as defined by each allelic profile.  For % mode

131     analyses, this is within a "group of interest or non-group of interest" framework. The outputs are

132     presented in the R console and optionally outputted to a tab-delimited format file. A facile method

133     to fully define the informative power of a SNP set derived by % analysis is to force the inclusion of

134     the SNPs into a *D* analysis, in which the user-defined SNP set size equals the number of included

135     SNPs, i.e., no additional SNPs are derived. This will reveal how the sequences assort in relation to

136     allelic profiles of the "forced included" SNPs. Alternatively, this can be done in reverse to assess the

137     performance of a *D* maximised SNP set to detect user-defined subsets of sequences with 100%

138     sensitivity. These functions provide considerable flexibility regarding the exploration of SNP sets.

## 6.1    Demonstration of MinSNPs functionality

140     To explore the potential utility of MinSNPs, we:

141      1. Determined the relationship between input alignment dimensions and the number and size

142        of output SNP sets, with running time;

143      2. Generated SNP sets of potential relevance to surveillance from orthologous SNP matrices

144        derived from genomic epidemiology studies in *Staphylococcus aureus* and *Plasmodium vivax*;

145      3. Explored the properties of SNPs identified with MinSNPs with respect to genome position

146        and relationship with coding sequences and;

147      4. Developed a pipeline to generate single orthologous SNP matrices from multiple short-read

148        data sets. This can support the analysis of large-scale comparative genome data using

149        MinSNPs.

### 6.1.1    Run-time Determinations

151 The relationships between the analysis time and dimensions of the input alignment, the number of

152 SNPs in the output SNP set, and the number of SNP sets in the output were determined. The

153 relationship was linear with respect to all three parameters. Examples of running time are shown in

154 **Error! Reference source not found.**. It was also shown that running MinSNPs using multiple cores

155    improves its performance. Complete data and code are shown in Supplementary Information 1

156    (https://figshare.com/s/696f696c232404f18a36). The faster run-time on a laptop as compared to a

157    high-performance cluster (HPC) was due to the simpler architecture of the machine; we note that

158    when the dimension of the alignments increases, the HPC's performance improves. So, given a

159    higher number of cores and increased memory available, a HPC can easily outperform a laptop.

160    *6.1.2*    **Derivation of SNP sets from a *Staphylococcus aureus* orthologous SNP matrix.**

161    To demonstrate MinSNPs' functionality, we analysed genome-wide orthologous SNP matrices to

162    identify 1. SNP sets diagnostic for a conserved lineage that is a potential surveillance target, 2. SNP

163    sets diagnostic for a broader phylogenetic lineage that encompasses the potential surveillance

164    target, and 3. SNP sets optimised with respect to $D$. For the latter, our interests were in the resolving

165    power (the $D$ value), and the concordance of the genotypes defined by the SNP sets with the

166    phylogeny indicated by the orthologous SNP matrix.

167    We first analysed a previously described orthologous SNP matrix composed of 20,651 SNPs from 162

168    *S. aureus* isolates, four *Staphylococcus argenteus* isolates, and *S. aureus* Mu50, which was the

169    reference genome for matrix construction (13). The isolates were from the STARRS study, which

170    revealed potential *S. aureus* transmission events involving haemodialysis patients, and potential

171    contacts in the clinical context environment, in the north of the Australian Northern Territory (13).

172    The STARRS study identified isolates of multilocus sequence typing (MLST) defined ST762 (clonal

173    complex (CC) 1), and were involved in transmission events leading to patient infections. ST762 is

174    vanishingly rare globally but was prevalent in the STARRs study. We therefore used the ST762

175    lineage identified in the STARRs study as a model for a potential surveillance target. Using MinSNPs

176    in % mode, we determined that 12 SNPs each individually discriminated all the ST762 isolates from

177    other isolates in the study, with 100% sensitivity and specificity (Supplementary Information 2

minSNP_paper_final

178    (https://figshare.com/s/746dd263140963185c53)). A BLAST analysis demonstrated that for each of

179    these SNPs, the alleles present in the ST762 isolates were not present in the public databases,

180    suggesting that these SNPs have generalised ability to discriminate ST762 from the remainder of the

181    *S. aureus* complex (Supplementary Information 3 (https://figshare.com/s/8adc2b14052ccb89dbed)).

182    The same procedure was used to derive SNP sets that discriminate the CC1 (ST1 and ST762) STARRS

183    isolates from the other isolates. It was found that there were 119 SNPs that each individually

184    provided 100% sensitivity and specificity (Supplementary Information 2). Similar to SNPs identified

185    for ST762, a BLAST analysis returned 61 specimens from Genbank; out of these 53 are CC1 with 3

186    false positives belonging to ST425 and 5 specimens untypeable by MLST.

187    We further used MinSNPs to derive 15 five-member SNP sets with maximised *D*. The *D* values

188    obtained ranged from 0.925 to 0.936, defining 16 to 21 genotypes. Concordance with phylogeny was

189    determined for two SNP sets (set 1 and 11) that were selected on the basis of having no SNPs in

190    common. Both SNP sets discriminate the major lineages defined by the STARRS SNP matrix (**Error!**

191    **Reference source not found.**, **Error! Reference source not found.**).

192    **6.1.3    Derivation of *Plasmodium vivax* SNP sets**
193    Given challenges associated with the large genome size and high proportions of 'contaminating'

194    human DNA, targeted SNP genotyping remains an important approach in *Plasmodium*

195    epidemiological tracking (16-18). MinSNPs was tested with a *P. vivax* orthologous SNP matrix

196    encompassing 259 isolates and 527,107 SNPs (19). The matrix encompassed heterozygote positions

197    (read as nucleotide ambiguities in MinSNPs) that enabled us to develop strategies to accommodate

198    this feature, which is common in polyclonal infections.

199    The data were generated from isolates collected from Malaysia, Thailand, and Indonesia, as part of a

200    study to identify changes in the *P. vivax* population as Sabah (Malaysia) approaches the elimination

201 of vivax malaria (19). In 183,509 of the SNPs, a nucleotide ambiguity code (where calls were

202 heterozygote) was assigned to at least one of these isolates.

203 As previously described, a subset of 26 specimens from Malaysia were near identical. These were

204 denoted "K2" strains reflecting isolates that were potentially undergoing clonal expansion (19). We

205 regarded these as a model surveillance target. SNPs that discriminated the K2 lineage were

206 identified with MinSNPs in % mode, with all the K2 specimens defined as the group of interest. All

207 183,509 positions where any of the sequences had an ambiguity code were excluded from the

208 analysis. The resulting analysis of 343,598 SNPs yielded 124 SNPs that each individually discriminated

209 the K2 lineage from all the other isolates in the matrix (Supplementary Information 4

210 https://figshare.com/s/db47a069aab93f3c615c)). Any of these 124 SNPs could potentially form the

211 basis of a K2 surveillance tool protocol, and using more than one of these SNPs may provide useful

212 redundancy to avoid false negatives due to undiscovered sequence diversity.

213 Next, SNPs that discriminated all Malaysian specimens from all other specimens were derived. To

214 streamline the analysis, only one K2 specimen was included. Also, three specimens that were

215 obtained in Malaysia but were likely to be imported from other regions based on their genomic

216 clustering patterns (PY0045-C, PY0004-C and PY0120-C) were omitted from the group of interest.

217 Initially, we confined the analysis to the 343,598 SNPs that do not encompass any ambiguity codes.

218 This was not successful. The maximum % obtained from five SNPs was 0.265, meaning that 73.5% of

219 the non-Malaysian specimens were not discriminated from the Malaysian specimens

220 (Supplementary Information 4). A different protocol was then adopted. Prior to MinSNPs analysis,

221 ambiguity codes were transformed into the major allele at that position (Supplementary Information

222 4 (https://figshare.com/s/db47a069aab93f3c615c)). Fortuitously, in all cases, the major allele was

223 consistent with the ambiguity code. After MinSNPs analysis, the relationship between the allelic

224 profiles and isolate was determined using the untransformed matrix. The untransformed matrix can

225  define allelic profiles that include ambiguity codes. Any specimens that had such an allelic profile,

226  i.e., they had an ambiguity code at a SNP within the SNP set being assessed, were classified as

227  untypeable by that SNP set. Typeability was therefore a criterion we used for assessing SNP sets,

228  although we do note that typeability is likely a function of specimen quality and/or whether the

229  specimen contained a mixture of strains. It is not an inherent property of a pure *P. vivax* clone.

230  This approach to identifying SNPs that discriminated Malaysian specimens was successful. Two sets

231  of two SNPs were identified, each of which discriminated all Malaysian specimens from all other

232  typable specimens. For one SNP set, 20 specimens (7.72%) were untypeable, and for the other, the

233  number of untypeable specimens was 22 (8.49%). All the Malaysian specimens were typable with

234  both SNP sets. The reason for the superior result from the matrix with ambiguity codes transformed

235  is unclear. However, we note that the MinSNPs' requirement in % mode that SNP sets provide 100%

236  sensitivity for the group of interest, is a stringent constraint. A false negative defined by a single

237  member of a group of interest disqualifies a position from inclusion in a SNP set. Being able to

238  capture more diversity for the analysis by using the transformation procedure also appears to have

239  been critical. A possible work-around for this constraint on SNP selection is to run separate analyses,

240  each with subsets of the group of interest.

241  We then used MinSNPs to derive *D* maximised SNP sets from the *P. vivax* alignment. Both the

242  approaches described above for accommodating ambiguity codes were used. Five SNP sets, each

243  comprising five SNPs were derived using each approach. When all the positions that encompassed at

244  least one ambiguity code were excluded from the analysis, the *D* values obtained were 0.751, 0.750,

245  0.572, and 0.564 (two sets). The most discriminatory SNP set (D = 0.751) was investigated further. It

246  was determined that the matrix defined eight allelic profiles. Although this number of profiles and

247  the *D* value do not indicate high discrimination, there was close concordance between allelic profile

248  and          country          of          origin          ((Supplementary          Information          4

249    (https://figshare.com/s/db47a069aab93f3c615c), **Error! Reference source not found.**). Thus, within

250    the context of the diversity defined by the input matrix, five SNPs can accurately reveal *P. vivax*

251    country of origin. When the analysis was repeated with the transformed ambiguity codes, very

252    different results were obtained. The *D* values were between 0.958 to 0.960, which is considerably

253    higher than in the previous experiment. Consistent with this, the SNP sets defined 31-32 allelic

254    profiles. The numbers of specimens defined as untypeable were significant, ranging from 64 to 68

255    (25%-26% of specimens). The concordance between country of origin was poor. Even with the larger

256    number of allelic profiles, there were numerous instances of specimens from different countries

257    having the same profile. A likely explanation is that positions that encompass ambiguity codes are

258    polymorphic within countries. Such SNPs are more likely to generate ambiguity codes because both

259    alleles may be present in a mixed infection. The exclusion of these positions will enrich for SNPs that

260    separate specimens from different countries and are monomorphic within countries. This would be

261    expected to facilitate the derivation of SNP sets that indicate the country of origin.

262    *6.1.4*    **Derivation of SNP sets from merged matrices.**

263    We further demonstrated the ability of MinSNPs to analyse large datasets. To this end, we obtained

264    additional *S. aureus* data collected through different initiatives (Bioprojects from Genbank:

265    PRJEB3174 (20, 21), PRJEB32286 (21), and PRJNA400143 (22)) and created a large orthologous SNP

266    matrix using a modification of the SPANDx pipeline (23) (Supplementary Information 5

267    (https://figshare.com/s/aadf860f3cfd9c416e3f)). The matrix encompasses 3,279 isolates (including

268    the reference genome Mu50) and 164,335 SNP positions. We then used this matrix to validate the

269    SNPs discriminating both ST762 and CC1 obtained earlier using only STARRS dataset. It was found

270    that apart from one SNP set, all the previously identified single SNP sets retained 100% sensitivity

271    and specificity for ST762 with this large data set.  However, two of the SNPs were not present in the

272    matrix. For CC1 (ST1, ST762, ST2851, ST2981), most of the previously identified SNP sets were not

273    fully present in the matrix (i.e., the STARRS derived sets often included positions that were not

274    included in the merged matrix due to quality filtering). For similar reasons, not all the members of

275    previously identified high-D SNPs-sets were present in the new matrix, and no meaningful

276    comparison between the previous analysis and current analysis could be made (see Supplementary

277    information 5 (https://figshare.com/s/aadf860f3cfd9c416e3f )).

278    We also reran the same tasks in 6.1.2 with the matrix. We identified 50 individual SNPs and 50 two-

279    member SNP sets that discriminate all ST762 isolates from all others. We similarly identified 39

280    individual SNPs and 61 two-member SNP sets (100 SNPs sets) that discriminate all CC1 isolates from

281    all others.

282    We then experimented with the $D$ mode analysis to accomplish two different tasks. First, we

283    attempted to identify SNPs that discriminated all CCs from each other. To accomplish this, all the

284    variant positions between isolates within the same CC were identified and recorded. A reduced

285    matrix was then constructed that contained only a single isolate from each of the CCs. We then

286    excluded from analysis all the previously recorded variant positions within CCs, before running a $D$

287    mode search. It was found that a minimum of seven SNPs were required to discriminate all 33 CCs

288    from each other. MinSNPs was tasked to provide 200 alternative SNP sets that achieved a $D$ of 1.0.

289    Of these, 165 of the sets had seven members; the remaining had eight members.

290    Next, we explored the resolving power of SNP sets identified simply to maximise $D$, without

291    reference to CC. Similarly, we identified five high-$D$ 10-SNP sets (Supplementary Information 5). Prior

292    to running MinSNPs analysis, all but a subset of 100 CC22 isolates were randomly selected to be

293    included in the input matrix to avoid overly biasing the analysis to include SNPs that discriminated

294    within CC22. We obtained SNP sets with $D$ values (recalculated using the entire matrix) ranging from

295    0.6314 to 0.6461. We selected the SNP set with the highest $D$ value and constructed the allelic

296    profile with the first 5 SNPs (see Supplementary Information 5

297    (https://figshare.com/s/aadf860f3cfd9c416e3f)). As expected from the similar experiment

298    performed with the smaller STARRS data set, there was close but imperfect correspondence

299    between CC and allelic profile, even though there was no reference to CC in the SNP derivation

300    procedure (see Supplementary Information 5 for comparison).

301    In summary, MinSNPs provides a flexible means for deriving SNP sets from sequence alignments that

302    are optimised for lineage-specific or generalised resolving power. We have demonstrated its utility

303    using large data sets, where one such data set was a SNP matrix assembled from multiple *S. aureus*

304    bioprojects, using a modified pipeline that we also report here. This provides the potential for

305    assembling matrices encompassing the global diversity of microorganisms and mining there for

306    optimised marker sets.

# 7.   Author statements

## 7.1    Authors and contributors

309    KSH: Data curation, Formal analysis, Investigation, Methodology, Software, Visualisation, Writing-

310    original draft, Writing – review and editing.

311    DCH: Data curation, Methodology, Resources, Project Administration, Supervision, Writing- review

312    and editing.

313    SA: Methodology, Resources, Supervision, Writing- review and editing.

314    PS: Funding acquisition, Methodology, Software, Supervision, Writing-Review and Editing

315    PMG: Conceptualisation, Formal analysis, Funding acquisition, Investigation, Methodology, Project

316    Administration, Resources, Supervision, Writing – original draft, Writing- review and editing.

317    **7.2    Conflicts of interest**

318    The authors declare there are no conflicts of interest.

319    **7.3    Funding information**

324    *7.4*    **Ethical approval**

325    This work does not involve human or animal research. The research team received written

326    confirmation of this from the Human Research Ethics Committee for the Northern Territory

327    Government Department of Health and the Menzies School of Health research.

328    *7.5*    **Consent for publication**

329    Not applicable

330    **7.6    Acknowledgements**

## 8.   *References*

335    1.        Kim T-W, Jang Y-H, Jeong MK, Seo Y, Park CH, Kang S, et al. Single-nucleotide polymorphism-

336    based epidemiological analysis of Korean Mycobacterium bovis isolates. Journal of Veterinary
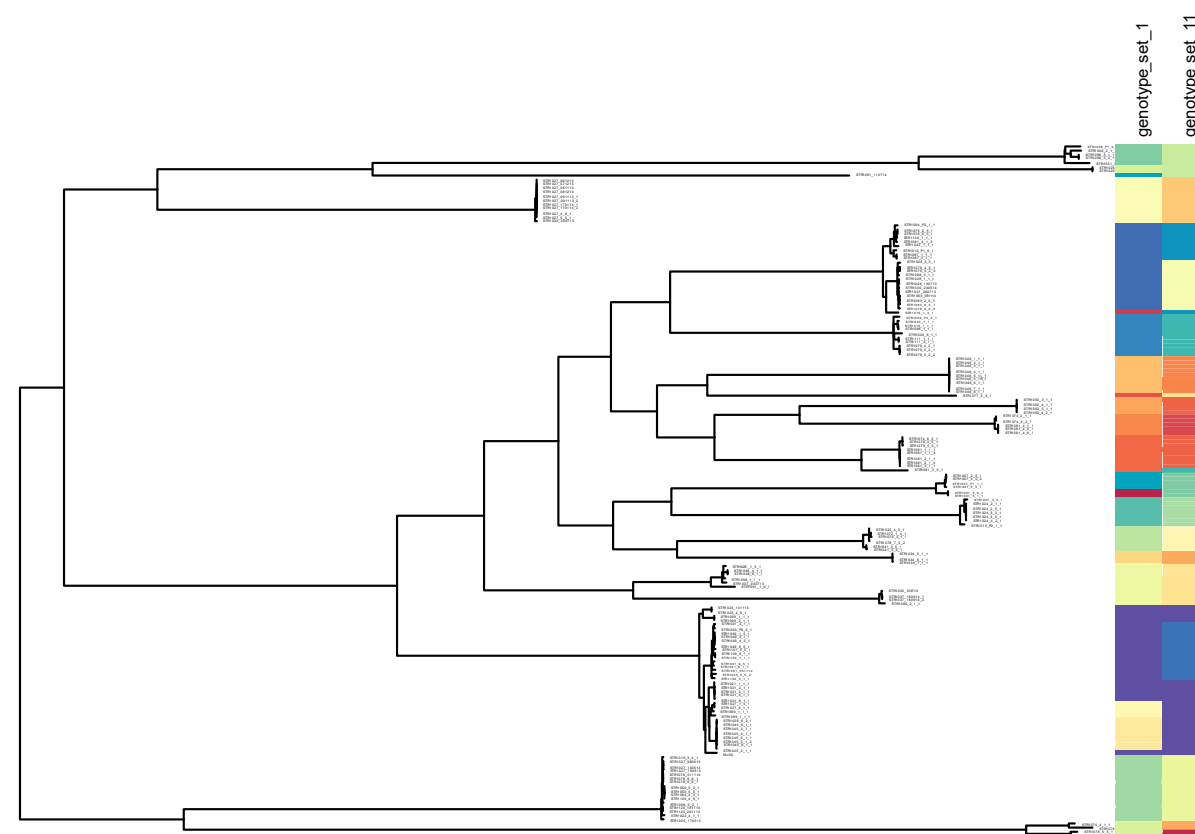
337    Science. 2021;22(2):e24-e.

338    2.    Napier G, Campino S, Merid Y, Abebe M, Woldeamanuel Y, Aseffa A, et al. Robust barcoding

339    and identification of Mycobacterium tuberculosis lineages for epidemiological and clinical studies.

340    Genome Medicine. 2020;12(1):114-.

341    3.    Vorimore F, Aaziz R, de Barbeyrac B, Peuchant O, Szymańska-Czerwińska M, Herrmann B, et

342    al. A New SNP-Based Genotyping Method for C. psittaci: Application to Field Samples for Quick

343    Identification. Microorganisms. 2021;9(3):625-.

344    4.    Jehanne Q, Pascoe B, Bénéjat L, Ducournau A, Buissonnière A, Mourkas E, et al. Genome-

345    Wide Identification of Host-Segregating Single-Nucleotide Polymorphisms for Source Attribution of

346    Clinical Campylobacter coli Isolates. Applied and Environmental Microbiology. 2020;86(24):e01787–

347    20-e–20.

348    5.    Kato CY, Chung IH, Robinson LK, Eremeeva ME, Dasch GA. Genetic typing of isolates of

349    Rickettsia typhi. PLoS Neglected Tropical Diseases. 2022;16(5):e0010354-e.

350    6.    Rahman M-M, Lim S-J, Park Y-C. Development of Single Nucleotide Polymorphism (SNP)-

351    Based Triplex PCR Marker for Serotype-specific Escherichia coli Detection. Pathogens.

352    2022;11(2):115-.

353    7.    Jacob CG, Thuy-Nhien N, Mayxay M, Maude RJ, Quang HH, Hongvanthong B, et al. Genetic

354    surveillance in the Greater Mekong subregion and South Asia to support malaria control and

355    elimination. eLife. 2021;10:e62997-e.

356    8.    Robertson GA, Thiruvenkataswamy V, Shilling H, Price EP, Huygens F, Henskens FA, et al.

357    Identification and interrogation of highly informative single nucleotide polymorphism sets defined

358    by bacterial multilocus sequence typing databases. J Med Microbiol. 2004;53(Pt 1):35-45.

359    9.    Price EP, Inman-Bamber J, Thiruvenkataswamy V, Huygens F, Giffard PM. Computer-aided

360    identification of polymorphism sets diagnostic for groups of bacterial and viral genetic variants. BMC

361    Bioinformatics. 2007;8:278.

362    10.    Tong SYC, Xie S, Richardson LJ, Ballard SA, Dakh F, Grabsch EA, et al. High-Resolution Melting

363    Genotyping of Enterococcus faecium Based on Multilocus Sequence Typing Derived Single

364    Nucleotide Polymorphisms. PLOS ONE. 2011;6(12):e29189-e.

365    11.    Price EP, Inman-Bamber J, Thiruvenkataswamy V, Huygens F, Giffard PM. Computer-aided

366    identification of polymorphism sets diagnostic for groups of bacterial and viral genetic variants. BMC

367    Bioinformatics. 2007;8(1):278-.

368    12.    Giffard PM, Andersson P, Wilson J, Buckley C, Lilliebridge R, Harris TM, et al. CtGEM typing:

369    Discrimination of Chlamydia trachomatis ocular and urogenital strains and major evolutionary

370    lineages by high resolution melting analysis of two amplified DNA fragments. PLOS ONE.

371    2018;13(4):e0195454-e.

372    13.    Holt DC, Harris TM, Hughes JT, Lilliebridge R, Croker D, Graham S, et al. Longitudinal whole-

373    genome based comparison of carriage and infection associated Staphylococcus aureus in northern

374    Australian dialysis clinics. 2021;16(2):e0245790-e.

375    14.    Lilliebridge RA, Tong SY, Giffard PM, Holt DC. MLST based Staphylococcus aureus typing

376    scheme using high-resolution melting analysis of SNP nucleated PCR fragments. The clinical and

377    molecular epidemiology of community-associated Staphylococcus aureus in northern Australia.

378    2010:119-.

379    15.    Robertson G, Thiruvenkataswamy V, Shilling H, Price E, Huygens F, Henskens F, et al.

380    Identification and Interrogation of Highly Informative Single Nuceotide Polymorphism Sets Defined

381    by Bacterial Multilocus Sequence Typing Databases. Journal of medical microbiology. 2004;53:35-45.

382    16.    Noviyanti R, Miotto O, Barry A, Marfurt J, Siegel S, Thuy-Nhien N, et al. Implementing

383    parasite genotyping into national surveillance frameworks: feedback from control programmes and

384    researchers in the Asia–Pacific region. BioMed Central; 2020.

385    17.    Fola AA, Kattenberg E, Razook Z, Lautu-Gumal D, Lee S, Mehra S, et al. SNP barcodes provide

386    higher resolution than microsatellite markers to measure Plasmodium vivax population genetics.

387    Malaria Journal. 2020;19:375-.

388    18.    Diez Benavente E, Campos M, Phelan J, Nolder D, Dombrowski JG, Marinho CRF, et al. A

389    molecular barcode to inform the geographical origin and transmission dynamics of Plasmodium vivax

390    malaria. PLoS Genetics. 2020;16(2):e1008576-e.

391    19.    Auburn S, Benavente ED, Miotto O, Pearson RD, Amato R, Grigg MJ, et al. Genomic analysis

392    of a pre-elimination Malaysian Plasmodium vivax population reveals selective pressures and

393    changing transmission dynamics. Nature Communications. 2018;9(1):2585-.

394    20.    Toleman MS, Reuter S, Coll F, Harrison EM, Blane B, Brown NM, et al. Systematic

395    Surveillance Detects Multiple Silent Introductions and Household Transmission of Methicillin-

396    Resistant Staphylococcus aureus USA300 in the East of England. The Journal of Infectious Diseases.

397    2016;214(3):447–53-–53.

398    21.    Coll F, Raven KE, Knight GM, Blane B, Harrison EM, Leek D, et al. Definition of a genetic

399    relatedness cutoff to exclude recent transmission of meticillin-resistant Staphylococcus aureus: a

400    genomic epidemiology analysis. The Lancet Microbe. 2020;1(8):e328–e35-e–e35.

401    22.    Manara S, Pasolli E, Dolce D, Ravenni N, Campana S, Armanini F, et al. Whole-genome

402    epidemiology, characterisation, and phylogenetic reconstruction of Staphylococcus aureus strains in

403    a paediatric hospital. Genome medicine. 2018;10(1):1–19-1–.

404    23.    Sarovich DS, Price EP. SPANDx: a genomics pipeline for comparative analysis of large haploid

405    whole genome re-sequencing datasets. BMC research notes. 2014;7(1):1–9-1–9.

## 9. Figures and tables

Figure 1: STARRS: Phylogeny and genotypes as defined by high-D SNP sets 1 and 11.



The phylogenetic tree was taken from the original paper (13) and labelled with two newly identify

high-*D* SNP sets. (https://microreact.org/project/minsnps-starrs). High-*D* SNP sets 1 and 11 include

positions 111760, 1925985, 2663300, 2683490, 124088, and 539419, 1413096, 1146945, 2184528,

1577370 of the Mu50 reference genome respectively.

minSNP_paper_final

**MICROBIOLOGY SOCIETY**

413    Table 1: Input alignment dimensions versus run time.

| Input alignment dimensions | Mode | Number of SNPs in SNP set | Running time HPC (s) | | Running time Laptop (s) | |
|---|---|---|---|---|---|---|
| | | | 2 Cores | 8 cores | 2 cores | 8 cores |
| 167 isolates; 20,651 SNPs | % | 1 | 31.926s | 20.907s | 16.809s | 7.027s |
| | D | 3 | 93.186s | 60.749s | 49.029s | 21.662s |
| | D | 5 | 157.831s | 105.136s | 85.098s | 35.363s |

414 Table 2: STARRS: Breakdown of CC/Singletons for genotypes defined by SNP sets 1 and 11. The distinction between singletons and CCs is somewhat

415 arbitrary. The CCs labelled with "*" were present only as the CC founder ST in the STARRS isolates. Column SA refers to *S. argenteus*.

416 **Table 2a Breakdown of CC/Singletons for genotypes defined by SNPs set 1**

| Genotype | SNPs set 1 (111760, 124088, 1925985, 2663300, 2683490) | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CC1 | CC5 | CC6 | CC8 | CC12* | CC15 | CC20* | ST30 | CC45 | CC72 | CC78 | CC93 | CC97 | CC101* | CC121* | ST834 | SA | Unknown |
| 1 | 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 24 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 16 | 0 | 0 | 0 | 0 | 0 | 0 |
| 18 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 |

| Genotype | CC1 | CC5 | CC6 | CC8 | CC12* | CC15 | CC20* | ST30 | CC45 | CC72 | CC78 | CC93 | CC97 | CC101* | CC121* | ST834 | SA | Unknown |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 19 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 0 | 0 |
| 20 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 11 | 0 | 0 | 0 |
| 21 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 0 | 0 |

417

**Table 2b Breakdown of CC/Singletons for genotypes defined by SNPs set 11**

| Genotype | SNPs set 11 (539419, 1146945, 1413096, 1577370, 2184528) | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | CC1 | CC5 | CC6 | CC8 | CC12* | CC15 | CC20* | ST30 | CC45 | CC72 | CC78 | CC93 | CC97 | CC101* | CC121* | ST834 | SA | Unknown |
| 1 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 2 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 3 | 0 | 22 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 4 | 0 | 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 5 | 0 | 0 | 5 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 6 | 0 | 0 | 0 | 10 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 7 | 0 | 0 | 0 | 0 | 9 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 8 | 0 | 0 | 0 | 0 | 0 | 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 9 | 0 | 0 | 0 | 0 | 0 | 0 | 4 | 0 | 0 | 8 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 10 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 |
| 11 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 7 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 13 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 16 | 0 | 0 | 0 | 0 | 0 | 0 |
| 14 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 6 | 0 | 0 | 0 | 0 | 0 |
| 15 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 7 | 0 | 0 | 0 | 0 |
| 16 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 11 | 0 | 0 | 0 |
| 17 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 3 | 2 | 0 |

minSNP_paper_final

419    Table 3: *P. vivax*: Genotypes defined by high-*D* SNP set 1 (ambiguity codes excluded vs substituted)

420    **Table 1a Genotypes defined by high-D SNPs set 1 (ambiguity code excluded)**

| Genotype | Excluded -(340505 [Chr 13], 460741 [Chr 12], 854772 [Chr 10], 531315 [Chr 6], 2100572 [Chr 12]) | | | |
|---|---|---|---|---|
| | Malaysia | Thailand | Indonesia | Imported |
| 1 | 26 | 0 | 0 | 0 |
| 2 | 17 | 1 | 0 | 0 |
| 3 | 3 | 3 | 0 | 0 |
| 4 | 1 | 91 | 0 | 1 |
| 5 | 0 | 9 | 0 | 0 |
| 6 | 1 | 0 | 80 | 2 |
| 7 | 0 | 0 | 11 | 0 |
| 8 | 0 | 0 | 9 | 0 |
| 9 | 0 | 0 | 3 | 0 |
| 10 | 0 | 0 | 1 | 0 |

421

422    **Table 3b Genotypes defined by high-D SNPs set 1 (ambiguity code substituted)**

| Genotype | Substituted (1269895 [Chr 14], 1240935 [Chr 13], 1812716 [Chr 11], 1717060 [Chr 9], 1141805 [Chr 10]) | | | |
|---|---|---|---|---|
| | Malaysia | Thailand | Indonesia | Imported |
| 1 | 26 | 0 | 0 | 0 |
| 2 | 5 | 0 | 1 | 0 |
| 3 | 3 | 5 | 0 | 0 |
| 4 | 3 | 0 | 2 | 1 |
| 5 | 2 | 0 | 5 | 0 |
| 6 | 1 | 7 | 0 | 0 |
| 7 | 1 | 5 | 0 | 0 |
| 8 | 1 | 0 | 4 | 0 |
| 9 | 0 | 8 | 0 | 0 |
| 10 | 0 | 8 | 0 | 0 |
| 11 | 0 | 7 | 0 | 0 |
| 12 | 0 | 6 | 0 | 0 |
| 13 | 0 | 5 | 0 | 1 |
| 14 | 0 | 5 | 0 | 0 |
| 15 | 0 | 5 | 0 | 0 |
| 16 | 0 | 5 | 0 | 0 |
| 17 | 0 | 3 | 0 | 0 |
| 18 | 0 | 3 | 0 | 0 |
| 19 | 0 | 2 | 0 | 0 |

| | | | | |
|---|---|---|---|---|
| 20 | 0 | 1 | 6 | 0 |
| 21 | 0 | 1 | 4 | 0 |
| 22 | 0 | 1 | 0 | 0 |
| 23 | 0 | 0 | 6 | 0 |
| 24 | 0 | 0 | 6 | 0 |
| 25 | 0 | 0 | 6 | 0 |
| 26 | 0 | 0 | 5 | 0 |
| 27 | 0 | 0 | 5 | 0 |
| 28 | 0 | 0 | 5 | 0 |
| 29 | 0 | 0 | 4 | 0 |
| 30 | 0 | 0 | 4 | 0 |
| 31 | 0 | 0 | 4 | 0 |
| 32 | 0 | 0 | 3 | 1 |

423

minSNP_paper_final