

1 **Global and regional dissemination and evolution of *Burkholderia pseudomallei***

2

3 Claire Chewapreecha<sup>1,2,3\*</sup>, Matthew T. G. Holden<sup>2,4</sup>, Minna Vehkala<sup>5</sup>, Niko  
4 Välimäki<sup>6</sup>, Zhirong Yang<sup>5</sup>, Simon R Harris<sup>2</sup>, Alison E. Mather<sup>7</sup>, Apichai Tuanyok<sup>8</sup>,  
5 Birgit De Smet<sup>9,10</sup>, Simon Le Hello<sup>11</sup>, Chantal Bizet<sup>12</sup>, Mark Mayo<sup>13</sup>, Vanaporn  
6 Wuthiekanun<sup>14</sup>, Direk Limmathurotsakul<sup>14,15,16</sup>, Rattanaphone Phetsouvanh<sup>17§</sup>, Brian  
7 G Spratt<sup>18</sup>, Jukka Corander<sup>5,19</sup>, Paul Keim<sup>20</sup>, Gordon Dougan<sup>1,2</sup>, David A. B.  
8 Dance<sup>16,17,21</sup>, Bart J Currie<sup>13</sup>, Julian Parkhill<sup>2</sup>, Sharon J. Peacock<sup>1,2,21\*</sup>

9

10 \* Correspondence should be addressed to Claire Chewapreecha

11 ([cchewapreecha@gmail.com](mailto:cchewapreecha@gmail.com)) and Sharon Peacock ([sharon.peacock@lshtm.ac.uk](mailto:sharon.peacock@lshtm.ac.uk))

12 § Deceased.

13

14 <sup>1</sup>Department of Medicine, University of Cambridge, UK

15 <sup>2</sup>Wellcome Trust Sanger Institute, Cambridge, UK

16 <sup>3</sup>Bioinformatics and Systems Biology Program, School of Bioresources and  
17 Technology, King Mongkut's University of Technology Thonburi, Thailand

18 <sup>4</sup>School of Medicine, University of St Andrew, UK

19 <sup>5</sup>Department of Mathematics and Statistics, University of Helsinki, Finland

20 <sup>6</sup>Department of Medical and Clinical Genetics, Genome-Scale Biology Research  
21 Program, University of Helsinki, Finland

22 <sup>7</sup>Department of Veterinary Medicine, University of Cambridge, UK

23 <sup>8</sup>Emerging Pathogens Institute, University of Florida, USA

24 <sup>9</sup>Department of Clinical Sciences, Institute of Tropical Medicine, Antwerp, Belgium

25 <sup>10</sup>Laboratory of Microbiology, Faculty of Sciences, Ghent University, Belgium

26 <sup>11</sup>Department of Infection and Epidemiology, Enteric bacteria pathogen Unit, Institut  
27 Pasteur, Paris, France

28 <sup>12</sup>Department of Microbiology, Collection of Institut Pasteur, Institut Pasteur, Paris,  
29 France

30 <sup>13</sup>Global and Tropical Health Division, Menzies School of Health Research, Charles  
31 Darwin University and Royal Darwin Hospital, Darwin, Australia

32 <sup>14</sup>Mahidol-Oxford Tropical Medicine Research Unit, Faculty of Tropical Medicine,  
33 Mahidol University, Bangkok, Thailand

34 <sup>15</sup>Department of Tropical Hygiene, Faculty of Tropical Medicine, Mahidol  
35 University, Bangkok, Thailand

36 <sup>16</sup>Centre for Tropical Medicine & Global Health, University of Oxford, UK

37 <sup>17</sup>Lao-Oxford-Mahosot Hospital-Wellcome Trust Research Unit, Microbiology  
38 Laboratory, Mahosot, Vientiane, Lao PDR

39 <sup>18</sup>Department of Infectious Disease Epidemiology, Imperial College, UK

40 <sup>19</sup>Department of Biostatistics, University of Oslo, Oslo, Norway

41 <sup>20</sup>Center for Microbial Genetics and Genomics, Northern Arizona University, USA

42 <sup>21</sup>London School of Hygiene and Tropical Medicine, UK

43 **The environmental bacterium *Burkholderia pseudomallei* causes an estimated**  
44 **165,000 cases of human melioidosis per year worldwide, and is also classified as a**  
45 **biothreat agent. We used whole genome sequences of 469 *B. pseudomallei* isolates**  
46 **from 30 countries collected over 79 years to explore its geographic transmission.**  
47 **Our data point to Australia as an early reservoir, with transmission to Southeast**  
48 **Asia followed by onward transmission to South Asia, and East Asia. Repeated**  
49 **reintroduction was observed within the Malay Peninsula, and between countries**  
50 **bordered by the Mekong river. Our data support an African origin of the**  
51 **Central and South American isolates with introduction of *B. pseudomallei* into**  
52 **the Americas between 1650 and 1850, providing a temporal link with the slave**  
53 **trade. We also identified geographically distinct genes/variants in Australasian**  
54 **or Southeast Asian isolates alone, with virulence-associated genes being among**  
55 **those overrepresented. This provides a potential explanation for clinical**  
56 **manifestations of melioidosis that are geographically restricted.**

57 *Burkholderia pseudomallei* is an environmental Gram-negative bacillus and  
58 the cause of melioidosis, a serious disease of humans and animals for which there is  
59 no licensed vaccine. Infection results from inoculation, ingestion or inhalation of *B.*  
60 *pseudomallei*, and is fatal in 10-40% of human cases<sup>1</sup>. To further understand the  
61 global dissemination of melioidosis, we sequenced 276 *B. pseudomallei* isolates  
62 cultured from humans with melioidosis or from the environment between 1935 and  
63 2013. These originated from 30 countries across Australasia, Asia, Africa and Central  
64 and South America. We added to this whole genome data available for a further 193  
65 *B. pseudomallei* isolates from Southeast Asia<sup>2</sup> and Australia<sup>3</sup>, giving a total dataset  
66 comprising 469 isolates (See Supplementary Data 1 for details of isolates and  
67 references). The genetic diversity of these isolates was captured by mapping short-  
68 read genome sequences against a core genome created from the two chromosomes of  
69 *B. pseudomallei* K96243<sup>4</sup>, and by extracting both core and accessory coding  
70 sequences from the assembled genomes (see methods). We employed three different  
71 approaches to outline the population structure: phylogenetic reconstructions using  
72 single nucleotide polymorphisms (SNPs) called from core genome mapping (Figure  
73 1a); SNPs from shared single-copy core genes (Supplementary Figure 1); and a tree-  
74 independent hierarchical Bayesian clustering (Supplementary Data 1).

75 All three approaches demonstrated a clear genetic distinction between isolates  
76 from Australasia and Asia (two areas where melioidosis is endemic), supporting  
77 previous findings<sup>5,6</sup>. Isolates from Australasia had longer phylogenetic branches  
78 compared to isolates from other regions, indicative of greater genetic diversity (Figure  
79 1a and Supplementary Figure 1). This was also observed from the pan-genome  
80 analysis<sup>7</sup>, which confirmed that the Australasian *B. pseudomallei* population had the  
81 highest rate of new gene discovery and the largest accessory genome (Figure 1b and

82 1c). Examination of data distribution confirmed that this finding was not related to  
83 different sampling periods or sequencing platforms used to generate the data  
84 (Supplementary Figure 2). These observations provide evidence for the hypothesis  
85 that Australia was an early reservoir for the current global *B. pseudomallei*  
86 population<sup>5,8</sup>, which is supported by the Australasian isolates being at the base of the  
87 tree (Supplementary Figure 1). An alternative explanation is that there have been  
88 repeated population bottlenecks outside Australia, but not within it. Figure 1a and  
89 Supplementary Figure 1 both delineated an apparent single transmission out of  
90 Australasia (consistent with previous findings<sup>9,10</sup>), and several independent  
91 transmission events from Southeast Asia to South Asia and East Asia. We also noted  
92 a monophyly and a single combined Bayesian cluster containing isolates from Africa  
93 and Central and South America, suggesting close ancestry (Figure 1a, Supplementary  
94 Figure 1 and Supplementary Data 1). The phylogenies also highlighted an African  
95 root for this group (100% bootstrap support), implying an African origin of the  
96 American isolates based on our sampling density.

97 We then estimated a timeline for the intercontinental and regional spread of *B.*  
98 *pseudomallei* by identifying and analysing 19 separate Bayesian clusters comprising  
99 isolates from Australia and Oceania (group 1), Asia (groups 2 to 18), and Africa and  
100 America (group 19). To improve our sensitivity to detect genetic variants, we  
101 remapped sequence reads from each cluster against a closely related reference  
102 genome (Supplementary Figure 3). After removing sequences that had been  
103 horizontally acquired by recombination<sup>11</sup>, temporal signals were determined for each  
104 cluster with the timeline estimated by BEAST<sup>12</sup> (Supplementary Figs 4, 5 and 6).  
105 Clock signals were captured for American isolates within the African-American  
106 cluster, and for four Asian clusters. The most recent common ancestor for the

107 American isolates was estimated to be 1806 or 1759 based on either chromosome I or  
108 II, respectively (combined 95% highest posterior density (HPD) interval of both  
109 chromosomes, 1682-1849) (Figure 2a). The introduction of *B. pseudomallei* into the  
110 Americas overlaps with the height of the slave trade between 1650 - 1850, during  
111 which an estimated 10-15 million people and related cargoes including  
112 environmentally contaminated food and water were transported from Africa to the  
113 Americas (Figure 2b)<sup>13,14</sup>. Dating of Asian clusters showed that recent common  
114 ancestors could be defined for three Malaysian-Singaporean clusters and one Thai –  
115 Laos cluster, all of which dated to the 20<sup>th</sup> century (Figure 2a). The most recent  
116 common ancestor of other Asian and Australasian clusters is very likely to pre-date  
117 these estimates, but dating of these deeper evolutionary events is less reliable.

118         Within the Asian isolates, the majority of Southeast Asian clusters either  
119 contained isolates from the Malay Peninsula (Malaysia and Singapore – here termed  
120 “the Malay sub-region”), or from countries bordered by the Mekong river (Thailand,  
121 Laos, Cambodia and Vietnam – here termed “the Mekong sub-region”)  
122 (Supplementary Figure 7a, Supplementary Data 1). To further examine this pattern,  
123 we estimated the number of times *B. pseudomallei* transitioned between Southeast  
124 Asian countries. This revealed a greater number of transitions within the same sub-  
125 regions than between sub-regions (two-tailed Mann-Whitney U test, p-value < 2.2x10<sup>-</sup>  
126 <sup>16</sup>) (Supplementary Figure 7b). The connectivity observed within sub-regions may be  
127 explained by geographical proximity, cultural links or trading networks associated  
128 with the Mekong river<sup>15,16</sup> (Figure 2c). In addition to an unequal number of  
129 transitions, *B. pseudomallei* may have spent different amounts of evolutionary time in  
130 these countries (total branch lengths of multiple sub-sampling phylogenetic trees)  
131 (Supplementary Figure 7c). Assuming a homogenous mutation rate, our results are

132 indicative of a higher proportion of evolutionary time spent in the Mekong versus the  
133 Malay sub-region (two-tailed Mann-Whitney U test,  $p$ -value  $< 2.2 \times 10^{-16}$ ), and  
134 possibly suggests that the Mekong sub-region has been a hotspot for *B. pseudomallei*  
135 evolution in the Southeast Asian endemic zone. It is possible that this observation  
136 may be influenced by evolutionary rate variation on each branch, but the local clock  
137 cannot be reliably assessed across this dataset.

138         The most common presentation of human melioidosis in both Asia and  
139 Australia is one or more of bacteremia, pneumonia and liver and/or splenic abscesses.  
140 By contrast, some of the less common clinical manifestations show geographical  
141 segregation, including encephalomyelitis in Australia. Moreover, mortality is lower in  
142 Australasia than Southeast Asia (10% versus 40%, respectively)<sup>17</sup>. Differences in  
143 human genetics and access to medical care including intensive care facilities are  
144 likely to contribute to different outcomes, but bacterial factors could also contribute to  
145 disease severity or to specific clinical manifestations. To investigate the genetic basis  
146 that might explain clinical differences between Australasia and Southeast Asia, we  
147 systematically screened for particular kmers (DNA words) that were enriched in  
148 Australasian isolates alone, or in Southeast Asian isolates alone using a kmer based  
149 GWAS<sup>18</sup> (see methods and Supplementary Data 2, 3 and 4). The strong link between  
150 the population structure and the geographical origin described above led us to omit  
151 population stratification in the GWAS analysis. Kmers were then clustered into loci  
152 based on their genetic proximity. This resulted in the identification of 468 and 14 loci  
153 that were specific to the Australasian and Southeast Asian population, respectively.  
154 Australasia- and Southeast Asia-specific loci were each distributed across multiple  
155 phylogenetic branches of their respective population (Supplementary Figure 8),  
156 suggesting that these were not solely driven by clonality in the population structure

157 but may have been independently acquired and/or lost on multiple occasions. The  
158 mechanisms that have driven these patterns will be the subject of further  
159 investigation.

160         Region-specific loci included those that may enhance survival and inter-  
161 bacterial competition in specific niches. They may also reflect virulence factors that  
162 contribute to the documented regionally distinct clinical manifestations. To facilitate  
163 the biological interpretation of these data, loci were categorised by the function of  
164 genes (COG), gene ontology (GO) and pathway terms. Some genes had no functional  
165 match in the curated database, but 64.3% could be assigned which revealed that  
166 region-specific genes were widely dispersed across multiple functions (Figure 3).  
167 Functional enrichment analyses highlighted elevated frequencies of the terms  
168 “secondary metabolite biosynthesis”, “translation”, “lipid transport and metabolism”  
169 and “defense mechanisms” among region-specific genes compared to random  
170 expectation from a reference genome (one-sided Fisher test p-value  $< 2.2 \times 10^{-16}$ ,  $<$   
171  $2.2 \times 10^{-16}$ ,  $1.86 \times 10^{-10}$  and  $9.07 \times 10^{-10}$  respectively, Supplementary Data 5). The  
172 latter contained several virulence genes involved in disease pathogenesis. Our results  
173 highlighted several virulence loci with known region-specific variations, including  
174 *Burkholderia thailandensis*-like flagellum and chemotaxis cluster (BTFC), and  
175 *Burkholderia mallei*-like *BimA* (*BmBimA*)<sup>19,20</sup>. Both BTFC and *BmBimA* facilitate  
176 bacterial motility inside host cells<sup>21,22</sup>, with the latter frequently detected in isolates  
177 associated with encephalomyelitis in Australia<sup>19</sup>. These findings validate our analytic  
178 approach and the ability to detect genetic variations based on geographical origin. The  
179 GWAS also identified unappreciated regional variations in well and less well  
180 characterised virulence loci (Supplementary Data 4), some examples of which are  
181 described below.



182 Filamentous hemagglutinin (*fha*) is a surface exposed and secreted protein that  
183 functions as an adhesin and immunomodulator across different bacterial species. In *B.*  
184 *pseudomallei*, the number of *fha* genes varies between isolates, and different  
185 combinations of *fha* genes have been observed between Australia and Thailand<sup>23</sup>.  
186 Furthermore, patients infected by *B. pseudomallei* with a specific *fha* variant are more  
187 likely to have infection associated with positive blood cultures<sup>19</sup>. We identified  
188 alternative adhesins/filamentous hemagglutinin variants in the Australasian  
189 population (Supplementary Data 4). For example, the BURPS668\_RS04895 variant in  
190 Australasian isolates differed from its non-Australasian ortholog by a group of kmers  
191 that clustered in an extended signal peptide for the Type V secretion system, and in  
192 hemagglutinin repeat domains (Supplementary Figure 9a). Such variation may alter  
193 protein secretion, binding affinity and specificity.

194 Intracellular pathogens have evolved various mechanisms for macrophage and  
195 immune evasion. Experimental evidence has shown that *B. pseudomallei* is capable of  
196 subverting antigen presentation and macrophage killing via polysaccharide capsule  
197 (CPS) and a type III secretion system (T3SS)<sup>24</sup>. We identified an Australasian-variant  
198 in CPS I (Supplementary Figure 9b), marked by kmers clustered in genes coding for  
199 two capsular polysaccharide export ABC transporter transmembrane proteins and  
200 putative sulfotransferase. We also identified variation in T3SS between the  
201 Australasian and Southeast Asian population (Supplementary Figure 9c). *B.*  
202 *pseudomallei* carries at least three clusters of T3SS, including T3SS-3 which is  
203 considered a virulence factor in mammalian infection. We noted genetic variants in  
204 T3SS-3 proteins *bsaU*, *bsaR*, *bsaP*, *bsaO*, an upstream region of a transcription factor  
205 *bprR* known to activate genes encoding structural components of T3SS-3<sup>25</sup>, and an  
206 oxygen-regulated invasion protein *orgA* in the Australasian population. Infection

207 assays using a macrophage cell line have shown reduced bacterial escape and lower  
208 intracellular bacterial survival of a *bsaU* mutant<sup>26</sup>, although the phenotype of  
209 geographical variants has not been established.

210 A distinctive feature of *B. pseudomallei* infection is the formation of multi-  
211 nucleated giant cells (MNGC), which results from cell membrane fusion between  
212 infected and uninfected host cells. This enables bacterial cell-to-cell spread while  
213 avoiding detection by host immunity. One of the key requirements for MNGC  
214 formation is a functional Type 6 secretion system cluster 1 (T6SS-1)<sup>24</sup>. We detected  
215 regional variation that extended from a known Australasian *BmBimA* variant to an  
216 upstream region of *virAG* regulator. This locus contains variations in hemolysin-  
217 coregulated protein (*hcp*), type VI secretion lysozyme-like protein (*tssE*), and ATP-  
218 dependent *clp* protease located on T6SS-1 (Supplementary Figure 9d). It remains to  
219 be seen whether region-specific variations in components of T6SS-1 and upstream of  
220 the *virA* regulator could affect disease pathogenesis.

221 In conclusion, our results indicate that movement of people and cargo has led  
222 to the dissemination of *B. pseudomallei*, a finding with implications for our  
223 increasingly globalised lifestyle. The carrier could have been contaminated soil, water  
224 or plants, or humans and other animals with clinical or sub-clinical disease. Given the  
225 frequency of *B. pseudomallei* transmission within Asia, it is striking that there appears  
226 to have been only one transmission event out of a diverse Australasian population into  
227 another geographical location. This might suggest that simple transmission is not  
228 sufficient, and that an adaptive bacterial event may also have been necessary. This  
229 could reflect the fact that the fauna of Australia and Southeast Asia are significantly  
230 different (the Wallace Line<sup>27</sup>). Identification of numerous bacterial genes or gene

231 variants that are geographically segregated provides a rich resource for biological  
232 studies of the basis for region-specific clinical syndromes in melioidosis.

233 **Methods**

234 **Bacterial collection and DNA sequencing.**

235 The global *B. pseudomallei* collection sequenced for this study contained 276 isolates  
236 from the environment and human disease. The rationale underpinning isolate selection  
237 from available global collections was to maximise distribution over time and  
238 geography, with representatives from each continent (see Supplementary Figure 2a).  
239 A very limited number of isolates had been stored and were available in areas where  
240 melioidosis is either uncommon or under-reported based on lack of microbiology  
241 infrastructure, which resulted in an unequal geographic representation. DNA libraries  
242 were prepared according to the Illumina protocol and sequenced on an Illumina  
243 HiSeq2000 or Miseq with paired-end runs to give a mean coverage of 84 reads per  
244 nucleotide (range 35 – 450). Publicly available sequence data for a further 193  
245 isolates (16 reference genomes, 76 Australasian isolates<sup>3</sup> and 101 Southeast Asian  
246 isolates<sup>2</sup>) and their accession numbers are also tabulated in Supplementary Data 1.

247

248 **Genome Assembly and Annotation.**

249 To control for potential contamination in each sample with other closely related  
250 species, taxonomic identity was assigned to all short reads and assemblies using  
251 Kraken<sup>28</sup>. Multilocus sequence typing (MLST) was derived from Illumina read data  
252 by mapping against the MLST sequence archive (<http://bpseudomallei.mlst.net/>).  
253 Unless previously assembled<sup>3</sup>, *de novo* assembly of short read data was performed  
254 using Velvet.<sup>29</sup> The kmer size was varied between 60% and 90% of the read length,  
255 and the assembly with the best N50 selected. Contigs shorter than the insert size  
256 length were filtered out. The sequence data were then used to further improve the  
257 assembly. Contigs were iteratively scaffolded using the process described in

258 Chewapreecha *et al.*<sup>30</sup>. As a QC step, reads were mapped back to the assembly using  
259 SMALT v. 0.7.4. (<http://www.sanger.ac.uk/resources/software/smalt/>). The  
260 assembly pipeline gave an average total length of 7,139,337 bp (range 6,744,467 –  
261 7,536,799) from 101 contigs (range 72 - 356) with an average contig length of 84,361  
262 bp (range 20,098– 192,188 bp) and an N50 of 223,075 (range 37,455 – 1,142,362).  
263 Gene predictions and annotations of draft reference genomes as well as other  
264 assemblies were performed using Prokka<sup>31</sup>. On average, 5,980 predicted coding  
265 sequences were assigned onto each genome (range 5,701 to 6,671 per each genome),  
266 falling within the similar range of a predicted 6,332 coding sequences in the first  
267 reference genome K96243 of 7.2 Mb.<sup>4,32</sup>

268

#### 269 **Pan-genome analysis.**

270 Based on annotated assemblies, a pan-genome was calculated for all 469 isolates  
271 using Roary<sup>7</sup>. An all-against-all comparison was performed using BLASTP and  
272 sequences clustered using a percentage identity of 92%, which was found to be a  
273 threshold that optimised specificity and sensitivity in this dataset (Supplementary  
274 Figure 10c and 10d). We identified a total of 25,812 predicted coding sequences  
275 (CDS), with 4,064 and 21,748 genes assigned to the core (present in 99% of isolates),  
276 and accessory (variably present) genome, respectively, which is comparable to that  
277 reported previously<sup>33</sup>. We used rarefaction curves to compare the number of predicted  
278 coding sequences as a function of the number of samples detected at different  
279 geographies (Figure 1b and Supplementary Figure 2c and 2d). A randomisation  
280 scheme with 1,000 permutations were employed to test our hypotheses about  
281 geographical diversity in gene contents. We also tested whether a greater rate of new  
282 gene discovery per number of samples sequenced in Australasia was biased by

283 different sampling timeframes or because the sequence data obtained from elsewhere  
284 were generated by different sequencing platforms. After sub-sampling the data  
285 (Supplementary Figure 2) to have equal representatives by year and sequencing  
286 quality, neither showed a change in the plot trajectory.

287

288 **Phylogeny based on shared single-copy core genes between *B. pseudomallei* and**  
289 ***B. thailandensis*.**

290 We repeated the pan-genome analysis described above with the inclusion of  
291 *Burkholderia thailandensis* genome E264 (accession numbers: NC\_007651.1 and  
292 NC\_007650.1), a closely related species that was used as an outgroup to root the tree.  
293 This demonstrated that 1,605 single-copy core genes were shared between *B.*  
294 *thailandensis* and *B. pseudomallei*. An approximate maximum likelihood  
295 phylogenetic tree was estimated by FastTree version 2.1.3<sup>34</sup> using GTR+CAT  
296 (General Time Reversible with per-site rate CATegories) model of approximation for  
297 site rate variation and was resampled 1,000 times (Supplementary Figure 1). The total  
298 number of single nucleotide polymorphic sites (SNPs) called was 127,421, of which  
299 69,473 SNPs (54.53%) represented differences between *B. thailandensis* and *B.*  
300 *pseudomallei*. This left 57,948 SNPs to resolve the *B. pseudomallei* population  
301 structure.

302

303 **Phylogeny based on core genome mapping of *B. pseudomallei***

304 A tree was constructed by mapping Illumina sequenced short reads to references  
305 using SMALT 0.7.4 (Figure 1a). Fully sequenced chromosomes and long reads  
306 sequenced by other platforms<sup>3</sup> were shredded to create 100 bp paired-end reads before  
307 mapping. Reads were mapped against the core genome of *B. pseudomallei* strain

308 K96243 (accession numbers BX571965 and BX571966) with bases called and  
309 aligned using a method previously described in Harris *et al.*<sup>35</sup> and Page *et al.*<sup>36</sup>.  
310 Genetic divergence compared with the K96243 core genome ranged from 0.73 to  
311 5.61%, and variants were identified at 324,637 SNPs (range 5,650 to 43,221 sites per  
312 isolate). A maximum-likelihood phylogeny was estimated with RAxML<sup>37</sup> using a  
313 general time reversible nucleotide substitution model with four gamma categories for  
314 rate heterogeneity and 100 bootstrap support.

315

### 316 **Hierarchical Bayesian clustering**

317 A tree-independent hierarchical Bayesian clustering with hierBAPS<sup>38,39</sup> was  
318 employed to determine the population structure generated from the core genome  
319 mapping alignment. This method allows the population to be sub-divided into groups  
320 with closely related genetic backgrounds and allows the recombination detection tool  
321 (Gubbins) to operate within its best performing range<sup>40</sup>. Except for the Australasian  
322 cluster (Group 1), which contained the highest amount of diversity for each isolate  
323 and could not be further sub-clustered, we continued the hierarchical clustering until  
324 the diversity observed in secondary or tertiary clusters fell within the limit of  
325 recombination detection (Supplementary Figure 10b). This resulted in 19 groups  
326 (Supplementary Data 1) for subsequent lineage-specific analyses. Except for Group  
327 15 and a bin cluster (35 isolates), Group 1 - 14 and 16 - 19 each formed a  
328 monophyletic group in the phylogeny (Figure 1a).

329

### 330 **Analysis of individual lineages.**

331 Evolutionary parameters and date of most recent common ancestors were determined  
332 for 19 clusters. For each cluster, closely related reference genomes were chosen for

333 mapping to increase variant calling sensitivity (Supplementary Figure 3). Where  
334 relevant reference genomes were not available as complete chromosomal contigs,  
335 draft reference genomes were created from *de novo* assemblies. One isolate within  
336 each of these clusters was selected, assembled and ordered relative to its closest  
337 reference using ABACAS v2.5.1<sup>41</sup> and ACT<sup>42</sup> followed by manual curation. Short  
338 reads from all members of each cluster were then mapped against this lineage-specific  
339 reference using SMALT 0.7.4. Bases were called and aligned with short insertions  
340 and deletions included using the method described in Harris *et al.*<sup>43</sup>. Recombination  
341 fragments were called and removed from the alignment using Gubbins<sup>11</sup>. A lineage-  
342 specific phylogeny was reconstructed using the remaining variants (Supplementary  
343 Figure 4).

344

#### 345 **Timeline reconstruction.**

346 We first tested for a positive correlation between date of isolation and root-to-tip  
347 distance obtained from a lineage-specific phylogeny with recombination removed  
348 using Path-O-Gen v1.4 (Supplementary Figure 5). Of 19 clusters, a consistent clock-  
349 like behaviour across both chromosomes was observed in a group of American  
350 isolates within the African-American cluster and five other Asian clusters (groups 4,  
351 5, 6, 7 and 8). Except for group 5 where the number of isolates were too low (n=4) to  
352 allow credible estimations, other clusters were analysed by BEAST v1.7<sup>12</sup> to  
353 determine the clock rate and the time when the most recent common ancestor  
354 emerged. We performed model selection on combinations of strict, relaxed log-  
355 normal, relaxed exponential, and random clock models and constant, exponential,  
356 logistic and skyline population models. For each, three independent chains were run  
357 for 50 million iterations, and sampled at every 1,000 generations. Models that failed



358 to converge based on visual inspection<sup>44</sup> of the trace files or had effective sampling  
359 size (ESS) values < 200 for key parameters were discarded. Stepping-stone and path-  
360 sampling analyses did not show appreciable differences between clock models,  
361 potentially suggesting that there may be insufficient rate variation within each group  
362 to warrant the use of a complex clock model. Thus, the strict clock with fewest  
363 parameters was employed to avoid over-fitting of parameters as suggested in <sup>45</sup>. We  
364 used the Bayesian skyline model as the tree prior to describing demographic history.  
365 Except for chromosome I of group 4 which did not achieve a credible ESS, the time  
366 calibrated phylogenetic trees, clock rates and time since most recent common ancestor  
367 (TMRCA) of estimated clusters are reported in Supplementary Figure 6.

368

369 Due to a small sample size used for each estimated cluster (American isolates within  
370 group 19: 9 isolates, group 4: 11 isolates, group 6: 24 isolates, group 7: 9 isolates, and  
371 group 8: 6 isolates), we also performed a date-randomised test as described in Murray  
372 *et al.*<sup>46</sup> to estimate the rigour of the true temporal signals compared to noise. For each  
373 tested cluster, we performed 1,000 permutations with the true date, but randomised  
374 root-to-tip distance. Regression coefficient  $R^2$  of the true data was ranked and  
375 compared to  $R^2$  of the randomised data (Supplementary Figure 5). Ranks of the true  
376 signals ranged from 34<sup>th</sup> (group 6 chromosome II) to 97<sup>th</sup> (group 8 chromosome II),  
377 suggesting that noise had an effect on a small dataset. Aside from small sample size,  
378 our clock rate on each chromosome for the clusters estimated by BEAST is consistent  
379 with previous estimates in *Burkholderia* species<sup>47</sup> and other bacteria<sup>35,48-50</sup>. This  
380 suggests that the results generated here are non-random.

381

382 **Ancestral state reconstruction on geographic locations of Southeast Asian**  
383 **isolates.**

384 Ancestral reconstruction was performed on the maximum likelihood global core  
385 genome phylogeny to assess the connectivity of isolates, and infer which population  
386 might act as source versus sink in Southeast Asia. To avoid sampling bias, we sub-  
387 sampled the phylogeny so that there were equal numbers of isolates from Thailand,  
388 Laos, Cambodia, Vietnam, Malaysia and Singapore (n=15 for each country), and  
389 resampled 1,000 times. Countries containing less than 15 isolates were excluded. We  
390 treated countries as discrete geographic characters. For each sub-sampled tree, we  
391 used stochastic character mapping *make.simmap* available in R package phytools  
392 v0.5-10<sup>51,52</sup> to estimate both the transitions between different geographical characters  
393 and the total time spent in each geographical character. Stochastic mapping was  
394 performed under an asymmetric model of character change for 1,000 simulations.

395

396 To assess the connectivity of isolates, we categorised geographical characters into two  
397 groups based on geographical proximity. The Mekong sub-region represents countries  
398 bordered by the Mekong river including Thailand, Laos, Cambodia and Vietnam; the  
399 Malay sub-region comprises Malaysia and Singapore. Changes between geographical  
400 characters were counted after grouping into two categories: 1) transitions within the  
401 same sub-region, and 2) transitions between sub-regions. The occurrence of  
402 transitions within and between the two sub-regions was compared using a two-tailed  
403 Mann-Whitney U test (Supplementary Figure 7b). To infer which population might  
404 act as the source, we compared the time spent in the Mekong and Malay sub-regions  
405 and compared this using a two-tailed Mann-Whitney U test (Supplementary Figure

406 7c). The choice of non-parametric Mann-Whitney U test over parametric test was due  
407 the violation of normally distributed data.

408

## 409 **Identification of distinct genes/variants in Australasian and Southeast Asian** 410 **populations.**

### 411 **Kmer-based GWAS without correction for population structure.**

412 We first considered the optimal approach to perform a GWAS for *B. pseudomallei*.  
413 Given the high level of genomic plasticity and large accessory genomes (Figure 1c),  
414 we concluded that a GWAS based on core genome SNPs as used elsewhere<sup>53,54</sup> would  
415 be sub-optimal as this fails to capture the extent of genetic variation. Instead, we used  
416 kmers (DNA words of length k) as an alternative to a SNP-based analysis. Unlike a  
417 traditional GWAS where genetic causes of particular phenotypes were identified  
418 while adjusting for population stratification, we employed GWAS to search for  
419 genetic markers in the Australasian and Southeast Asian populations, some of which  
420 may intrinsically define population structure. A control for population structure was  
421 thus omitted. Two independent GWAS runs were performed to search for variable  
422 kmers in the Australasia population alone (Australasia GWAS), and the Southeast  
423 Asian population alone (SEA GWAS). For both GWAS runs, the data were randomly  
424 divided into a discovery and a validation dataset. The Australasia GWAS comprised a  
425 set of 80 Australasia and 200 non-Australasian isolates, and was validated with 57  
426 Australasian and 132 non-Australasian isolates. Similarly, the SEA GWAS comprised  
427 a random set of 180 Southeast Asian and 105 non-Southeast Asian isolates, and was  
428 confirmed using 114 Southeast Asian and 65 non-Southeast Asian isolates. We used  
429 the reference-independent GWAS pipeline Seer by Lees *et al.*<sup>18</sup> to search for kmers  
430 with region-specific patterns. All kmers of length 9-100 bp were scanned from all

431 assembled reads using fsm-lite (<https://github.com/nvalimak/fsm-lite>). Only kmers  
432 seen in 5-95% of the total population were retained to reduce false positives from  
433 testing underpowered kmers. Seer<sup>18</sup> was performed on the discovery data using  
434 geographical origin of isolates (Australasia/ non-Australasia or Southeast Asia/ non-  
435 Southeast Asia) as binary phenotype ( $y$ ) and the presence/absence of each kmer as  
436 tested genotype  $X$ :

$$437 \quad \log\left(\frac{y}{1-y}\right) = X\beta$$

438 The direction of association (positive or negative) is described by  $\beta$ . Kmers with a  
439 conservative cut-off p-value  $< 10^{-8}$  in the logistic regression were considered further  
440 as suggested in Lees *et al.*<sup>18</sup>. Australasia and SEA GWAS yielded 77,787 and 43,663  
441 kmers, respectively, that were positively or negatively associated with Australasia or  
442 Southeast Asia populations (Supplementary Data 2). Among these, 42,521 kmers that  
443 were positively associated with Australasia were negatively associated with Southeast  
444 Asia (Supplementary Figure 11). Kmers that reached significance in the discovery  
445 data were confirmed in the validation data. To aid visualisation, the frequencies of  
446 5,000 randomly chosen kmers from the Australasia and SEA GWAS in the validation  
447 data have been plotted in Supplementary Figure 11.

448

#### 449 **Mapping and kmer clustering.**

450 Significant kmers were searched for an exact match in *de novo* assemblies and fully  
451 sequenced chromosomes using BLAT v. 34<sup>55</sup> with minimum match and score  
452 adjusted to cater for low complexity kmers as below.

453 `blat assembly kmers.query -minMatch=1 -minScore=10 output`

454 To facilitate biological interpretation, kmers were grouped into clusters based on their  
455 genetic distance. We defined the size of operons based on the length of transcription

456 fragments reported in Ooi *et al.*<sup>56</sup>. Any kmers located within 7.68 kb (the size of an  
457 operon covering 95<sup>th</sup> percentile of transcription fragments) were grouped together into  
458 a locus. On average, each locus had a median of 66 kmers (range 2 -11,072 kmers),  
459 with the size of the loci ranging from 40 – 70,684 bp. The binary patterns in size of  
460 region-specific loci (Supplementary Figure 12) likely reflect different scales of  
461 variation, with smaller and larger peaks corresponding to small-scale differences  
462 (including SNPs) and large-scale differences (including regions of mobile genetic  
463 elements incorporated via homologous recombination or site specific recombination),  
464 respectively. As the GWAS was not corrected for population structure, we further  
465 tested whether the predicted loci were subjected to clonality. The presence and  
466 absence of each locus (measured by % of detected kmers) were plotted against the  
467 phylogeny. Their scattering patterns across multiple branches suggested that region-  
468 specific loci were not strictly driven by a clonal population structure (Supplementary  
469 Figure 8).

470

#### 471 **COG, GO and pathway terms found in region-specific loci.**

472 We annotated the biological properties of kmers within coding regions using  
473 information from the functional categories (COG term), Gene Ontology (GO term),  
474 and pathway data (KEGG, InterPro and UniPathway), available from the  
475 *Burkholderia* Genome Database<sup>57</sup>. The reference genome Bp668 contained 40,986 out  
476 of 78,929 region-specific kmers, of which 23,565 overlapped with coding regions.  
477 One-sided Fisher's exact test was used to search for COG, GO and pathway terms in  
478 kmers that showed significant departure from random expectation in the Bp668  
479 genome. We tested kmers enrichment in 22 COG terms, 1,485 GO terms, and 408  
480 pathway terms using a strict Bonferroni correction with a required p-value of

481 0.01/1,915 =  $5.22 \times 10^{-6}$ . Significant COG terms were marked in Figure 3. Additional  
482 GO and pathway enrichment analyses are discussed in the supplementary note. Terms  
483 with significant deviation are tabulated in Supplementary Data 5.

484

#### 485 **Statistics and visualisation.**

486 Visualisation of phylogenetic trees and statistical analyses were performed in R<sup>58</sup>,  
487 iTOL<sup>59</sup>, and FigTree v 1.4.2 (<http://tree.bio.ed.ac.uk/software/figtree/>).

488

#### 489 **Data availability.**

490 New sequence data for the study isolates have been deposited in the ENA under study  
491 accession number ERP001193 and ERP002658, with the accession numbers for  
492 individual isolates listed in Supplementary Data 1. Supplementary Data 2-5 provide  
493 information that supports the data presented.

494

#### 495 **References**

496 Note: number 1 - 27 are in text references

497 1 Limmathurotsakul, D. *et al.* Predicted global distribution of *Burkholderia*  
498 *pseudomallei* and burden of melioidosis. *Nat Microbiol* 1,  
499 doi:10.1038/nmicrobiol.2015.8 (2016).

500 2 Nandi, T. *et al.* *Burkholderia pseudomallei* sequencing identifies genomic  
501 clades with distinct recombination, accessory, and epigenetic profiles. *Genome Res*  
502 25, 608 (2015).

503 3 Johnson, S. L. *et al.* Whole-Genome Sequences of 80 Environmental and  
504 Clinical Isolates of *Burkholderia pseudomallei*. *Genome Announc* 3,  
505 doi:10.1128/genomeA.01282-14 (2015).

506 4 Holden, M. T. *et al.* Genomic plasticity of the causative agent of melioidosis,  
507 *Burkholderia pseudomallei*. *Proc Natl Acad Sci U S A* 101, 14240-14245,  
508 doi:10.1073/pnas.0403302101 (2004).

509 5 Price, E. P. *et al.* Large-scale comparative genomics identifies unprecedented  
510 melioidosis cases in northern Australia caused by an Asian *Burkholderia*  
511 *pseudomallei* strain. *Appl Environ Microbiol*, doi:10.1128/AEM.03013-15 (2015).

512 6 Pearson, T. *et al.* Phylogeographic reconstruction of a bacterial species with  
513 high levels of lateral gene transfer. *BMC Biol* 7, 78, doi:10.1186/1741-7007-7-78  
514 (2009).

515 7 Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.  
516 *Bioinformatics* 31, 3691-3693, doi:10.1093/bioinformatics/btv421 (2015).

517 8 Dale, J. *et al.* Epidemiological tracking and population assignment of the non-  
518 clonal bacterium, *Burkholderia pseudomallei*. *PLoS Negl Trop Dis* 5, e1381,  
519 doi:10.1371/journal.pntd.0001381 (2011).

520 9 Gee, J. E., Allender, C. J., Tuanyok, A., Elrod, M. G. & Hoffmaster, A. R.  
521 *Burkholderia pseudomallei* type G in Western Hemisphere. *Emerg Infect Dis* 20, 682-  
522 684, doi:10.3201/eid2004.130960 (2014).

523 10 Sarovich, D. S. *et al.* Phylogenomic Analysis Reveals an Asian Origin for  
524 African *Burkholderia pseudomallei* and Further Supports Melioidosis Endemicity in  
525 Africa. *mSphere* 1, doi:10.1128/mSphere.00089-15 (2016).

526 11 Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of  
527 recombinant bacterial whole genome sequences using Gubbins. *Nucleic Acids Res* 43,  
528 e15, doi:10.1093/nar/gku1196 (2015).

529 12 Drummond, A. J., Suchard, M. A., Xie, D. & Rambaut, A. Bayesian  
530 phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol* 29, 1969-1973,  
531 doi:10.1093/molbev/mss075 (2012).

532 13 Kolchin, P. *American Slavery: 1619-1877*. (Penguin, 1995).

533 14 Thomas, H. *The Slave Trade, The Story of the Atlantic Slave Trade:1440-*  
534 *1870*. (Simon & Schuster Paperbacks, 1997).

535 15 Nguyen, T. D. *The Mekong River and the Struggle for Indochina: Water, War*  
536 *and Peace*. (Praeger Publishers, 1999).

537 16 Liu, J. H., Lawrence, B. & Ward, C. Social representations of history in  
538 Malaysia and Singapore: On the relationship between national and ethnic identity.  
539 *Asian J Soc Psychol* 5, 3-20 (2002).

540 17 Currie, B. J. Melioidosis: evolving concepts in epidemiology, pathogenesis,  
541 and treatment. *Semin Respir Crit Care Med* 36, 111-125, doi:10.1055/s-0034-  
542 1398389 (2015).

543 18 Lees, J. A. *et al*. Sequence element enrichment analysis to determine the  
544 genetic basis of bacterial phenotypes. *Nat Commun.* 7:12797  
545 doi:10.1038/ncomms12797 (2016).

546 19 Sarovich, D. S. *et al*. Variable virulence factors in *Burkholderia pseudomallei*  
547 (melioidosis) associated with human disease. *PLoS One* 9, e91682,  
548 doi:10.1371/journal.pone.0091682 (2014).

549 20 Tuanyok, A. *et al*. A horizontal gene transfer event defines two distinct groups  
550 within *Burkholderia pseudomallei* that have dissimilar geographic distributions. *J*  
551 *Bacteriol* 189, 9044-9049, doi:10.1128/JB.01264-07 (2007).



552 21 French, C. T. *et al.* Dissection of the *Burkholderia* intracellular life cycle  
553 using a photothermal nanoblade. *Proc Natl Acad Sci U S A* 108, 12095-12100,  
554 doi:10.1073/pnas.1107183108 (2011).

555 22 Benanti, E. L., Nguyen, C. M. & Welch, M. D. Virulent *Burkholderia* species  
556 mimic host actin polymerases to drive actin-based motility. *Cell* 161, 348-360,  
557 doi:10.1016/j.cell.2015.02.044 (2015).

558 23 Tuanyok, A. *et al.* Genomic islands from five strains of *Burkholderia*  
559 *pseudomallei*. *BMC Genomics* 9, 566, doi:10.1186/1471-2164-9-566 (2008).

560 24 Willcocks, S. J., Denman, C. C., Atkins, H. S. & Wren, B. W. Intracellular  
561 replication of the well-armed pathogen *Burkholderia pseudomallei*. *Curr Opin*  
562 *Microbiol* 29, 94-103, doi:10.1016/j.mib.2015.11.007 (2016).

563 25 Chen, Y. *et al.* Characterization and analysis of the *Burkholderia pseudomallei*  
564 *BsaN* virulence regulon. *BMC Microbiol* 14, 206, doi:10.1186/s12866-014-0206-6  
565 (2014).

566 26 Bast, A. *et al.* Caspase-1-dependent and -independent cell death pathways in  
567 *Burkholderia pseudomallei* infection of macrophages. *PLoS Pathog* 10, e1003986,  
568 doi:10.1371/journal.ppat.1003986 (2014).

569 27 Wallace, A. R. On the Physical Geography of the Malay Archipelago. *Journal*  
570 *of the Royal Geographical Society of London* 7, 205-212 (1863).

571 28 Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence  
572 classification using exact alignments. *Genome Biol* 15, R46, doi:10.1186/gb-  
573 2014-15-3-r46 (2014).

574 29 Zerbino, D. R. & Birney, E. Velvet: algorithms for de novo short read  
575 assembly using de Bruijn graphs. *Genome Res* 18, 821-829,  
576 doi:10.1101/gr.074492.107 (2008).

577 30 Chewapreecha, C. et al. Dense genomic sampling identifies highways of  
578 pneumococcal recombination. *Nat Genet* 46, 305-309, doi:10.1038/ng.2895  
579 (2014).

580 31 Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* 30,  
581 2068-2069, doi:10.1093/bioinformatics/btu153 (2014).

582 32 Nandi, T. et al. A genomic survey of positive selection in *Burkholderia*  
583 *pseudomallei* provides insights into the evolution of accidental virulence.  
584 *PLoS Pathog* 6, e1000845, doi:10.1371/journal.ppat.1000845 (2010).

585 33 Spring-Pearson, S. M. et al. Pangenome Analysis of *Burkholderia*  
586 *pseudomallei*: Genome Evolution Preserves Gene Order despite High  
587 Recombination Rates. *PLoS One* 10, e0140274,  
588 doi:10.1371/journal.pone.0140274 (2015).

589 34 Price, M. N., Dehal, P.S. & Arkin, A.P. FastTree 2 – Approximately  
590 Maximum-Likelihood Trees for Large Alignments. *PLoS One*, 5(3):e9490,  
591 doi:10.1371/journal.pone.0009490 (2010).

592 35 Harris, S. R. et al. Evolution of MRSA during hospital transmission and  
593 intercontinental spread. *Science* 327, 469-474, doi:10.1126/science.1182395  
594 (2010).

595 36 Page, A. J. et al. SNP-sites: rapid efficient extraction of SNPs from  
596 multiFASTA  
597 alignments. *Microbial Genomics* 2, doi:10.1099/mgen.0.000056  
598 (2016).

599 37 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and  
600 postanalysis  
601 of large phylogenies. *Bioinformatics* 30, 1312-1313,

602 doi:10.1093/bioinformatics/btu033 (2014).

603 38 Corander, J., Marttinen, P., Siren, J. & Tang, J. Enhanced Bayesian modelling  
604 in BAPS software for learning genetic structures of populations. *BMC*  
605 *Bioinformatics* 9, 539, doi:10.1186/1471-2105-9-539 (2008).

606 39 Cheng, L., Connor, T. R., Siren, J., Aanensen, D. M. & Corander, J.  
607 Hierarchical and Spatially Explicit Clustering of DNA Sequences with BAPS  
608 Software. *Mol Biol Evol* 30, 1224-1228, doi:10.1093/molbev/mst028 (2013).

609 40 Croucher, N. J. et al. Population genomics of post-vaccine changes in  
610 pneumococcal epidemiology. *Nat Genet* 45, 656-663, doi:10.1038/ng.2625  
611 (2013).

612 41 Assefa, S., Keane, T. M., Otto, T. D., Newbold, C. & Berriman, M. ABACAS:  
613 algorithm-based automatic contiguation of assembled sequences.  
614 *Bioinformatics* 25, 1968-1969, doi:10.1093/bioinformatics/btp347 (2009).

615 42 Carver, T. J. et al. ACT: the Artemis Comparison Tool. *Bioinformatics* 21,  
616 3422-3423, doi:10.1093/bioinformatics/bti553 (2005).

617 43 Harris, S. R. et al. Genome specialization and decay of the strangles pathogen,  
618 *Streptococcus equi*, is driven by persistent infection. *Genome Res* 25, 1360-  
619 1371, doi:10.1101/gr.189803.115 (2015).

620 44 Rambaut, A., Suchard, M. A., Xie, D. & Drummond, A. J. Tracer v1.6,  
621 <<http://beast.bio.ed.ac.uk/Tracer>> (2014).

622 45 Ho, S. Y. & Duchene, S. Molecular-clock methods for estimating evolutionary  
623 rates and timescales. *Mol Ecol* 23, 5947-5965, doi:10.1111/mec.12953 (2014).

624 46 Murray, G. G. R. et al. The effect of genetic structure on molecular dating and  
625 tests for temporal signal. *Methods in Ecology and Evolution*, 7, 80-89,  
626 doi:doi: 10.1111/2041-210X.12466 (2016).

627 47 Lieberman, T. D. et al. Parallel bacterial evolution within multiple patients  
628 identifies candidate pathogenicity genes. *Nat Genet* 43, 1275-1280,  
629 doi:10.1038/ng.997 (2011).

630 48 Mathers, A. J. et al. *Klebsiella pneumoniae* carbapenemase (KPC)-producing  
631 *K. pneumoniae* at a single institution: insights into endemicity from wholegenome  
632 sequencing. *Antimicrob Agents Chemother* 59, 1656-1663,  
633 doi:10.1128/AAC.04292-14 (2015).

634 49 Young, B. C. et al. Evolutionary dynamics of *Staphylococcus aureus* during  
635 progression from carriage to disease. *Proc Natl Acad Sci U S A* 109, 4550-  
636 4555, doi:10.1073/pnas.1113219109 (2012).

637 50 Croucher, N. J. et al. Rapid pneumococcal evolution in response to clinical  
638 interventions. *Science* 331, 430-434, doi:10.1126/science.1198545 (2011).

639 51 Revell, L. J. phytools: An R package for phylogenetic comparative biology  
640 (and other things). *Methods Ecol Evol* 3, 217-223 (2012).

641 52 Bollback, J. P. SIMMAP: stochastic character mapping of discrete traits on  
642 phylogenies. *BMC Bioinformatics* 7, 88, doi:10.1186/1471-2105-7-88 (2006).

643 53 Laabei, M. et al. Predicting the virulence of MRSA from its genome sequence.  
644 *Genome Res* 24, 839-849, doi:10.1101/gr.165415.113 (2014).

645 54 Chewapreecha, C. et al. Comprehensive identification of single nucleotide  
646 polymorphisms associated with beta-lactam resistance within pneumococcal  
647 mosaic genes. *PLoS Genet* 10, e1004547, doi:10.1371/journal.pgen.1004547  
648 (2014).

649 55 Kent, W. J. BLAT--the BLAST-like alignment tool. *Genome Res* 12, 656-664,  
650 doi:10.1101/gr.229202. (2002).

651 56 Ooi, W. F. et al. The condition-dependent transcriptional landscape of

652 *Burkholderia pseudomallei*. *PLoS Genet* 9, e1003795,  
653 doi:10.1371/journal.pgen.1003795 (2013).  
654 57 Winsor, G. L. et al. The *Burkholderia* Genome Database: facilitating flexible  
655 queries and comparative analyses. *Bioinformatics* 24, 2803-2804,  
656 doi:10.1093/bioinformatics/btn524 (2008).  
657 58 R Development Core Team R: A language and environment for  
658 statistical computing. R Foundation for Statistical Computing,  
659 Vienna, Austria. ISBN 3-900051-07-0, URL <http://www.R-project.org>. (2008)  
660 59 Letunic, I. & Bork, P. Interactive tree of life (iTOL) v3: an online tool for the  
661 display and annotation of phylogenetic and other trees. *Nucleic Acids Res* 44,  
662 W242-245, doi:10.1093/nar/gkw290 (2016).

663

#### 664 **Acknowledgements**

665 The authors thank the Wellcome Trust Sanger Institute library construction, sequence  
666 and core informatics teams, and Elizabeth Blane for their technical support. We thank  
667 Drs Tanistha Nandi and Patrick Tan at Genome Institute of Singapore; and Drs Erin  
668 Price and Derek Sarovich at Menzies School of Health Research, Australia for  
669 providing access to publically available WGS data. We thank the following people  
670 who provided isolates or DNA: Professor Nicholas Day, MORU, Faculty of Tropical  
671 Medicine, Mahidol University; Drs Paul Newton, Manivanh Vongsouvath, Mayfong  
672 Mayway, Viengmon Davong, Olay Lattana, Catrin Moore, Sayaphet Rattanavong and  
673 the directors and staff of Mahosot Hospital, Vientiane, Lao PDR; Dr Varun Kumar,  
674 Ankor Hospital for Children, Siem Reap, Cambodia; Dr James Campbell, Oxford  
675 University Clinical Research Unit, Ho Chi Minh City, Vietnam; Dr Hui Suk Wai,  
676 Ocean Park Corporation, Hong Kong SAR, China; Mr Chun Kham and Dr Thong

677 Phe, Sihanouk Hospital Centre of Hope, Phnom Penh, Cambodia; Dr Joost W.  
678 Wiersinga, Academic Medical Center (AMC), Amsterdam, the Netherlands; Professor  
679 Jan Jacobs, ITM, Antwerp, Belgium; Dr Julie E. Russell, National Collection of Type  
680 Cultures, UK; Dr Ty Pitt, NHS Blood and Transplant, UK; Mr Daniel Godoy,  
681 Imperial College, UK; Dr Stephane Emonet, Geneva University Hospitals,  
682 Switzerland; Dr Susan Morpeth, Middlemore Hospital, New Zealand; and Dr Jay  
683 Gee, CDC, USA. C.C. is a Sir Henry Wellcome post-doctoral Fellow (grant ref:  
684 107376/Z/15/Z). J.C., M.V., Z.Y. were supported by the COIN Centre of Excellence  
685 and Z.Y. by a HIIT post-doctoral fellowship. A.E.M. is supported by Biotechnology  
686 and Biological Sciences Research Council grant BB/M014088/1. B.G.S. was  
687 supported by the Wellcome Trust grant WT089472. D.A.B.D and R.P. are supported  
688 by the Wellcome Trust grants 106698/Z/14 and B9R00760. D.L. and V.W. are  
689 supported by the Wellcome Trust grant 089275/Z/09/Z. M.M. and B.J.C are  
690 supported by the Australian National Health and Medical Research Council through  
691 project grants #1046812 and #1098337. This publication presents independent  
692 research supported by the Health Innovation Challenge Fund (WT098600, HICF-T5-  
693 342), a parallel funding partnership between the Department of Health and Wellcome  
694 Trust. The views expressed in this publication are those of the author(s) and not  
695 necessarily those of the Department of Health or Wellcome Trust. This project was  
696 also funded by a grant awarded to the Wellcome Trust Sanger Institute (098051).

697

#### 698 **Author contributions**

699 A.T., B.D.S., S.L.H., C.B., M.M., V.W., D.L., R.P., B.G.S., P.K., D.A.B.D. and  
700 B.J.C. collected and provided the samples for the study. C.C. designed and performed  
701 the analyses. M.T.G.H, S.R.H., A.E.M., J.C., J.P. and G.D. designed and contributed

702 materials and analysis tools. M.V., N.V., Z.Y., and J.C. performed the kmer based  
703 analyses in the first draft. C.C. performed the kmer based analysis in the revised draft.  
704 Z.Y. and J.C. performed cluster analyses. S.J.P. was responsible for management of  
705 the study. S.J.P. and C.C. wrote the paper with input from all authors. All authors  
706 approved the manuscript prior to submission.

707

## 708 **Figure legends**

### 709 **Figure 1 The phylogeny and pan-genome of *B. pseudomallei***

710 Differences in level of bacterial diversity across different geographical origins:  
711 Australasia (green), Asia (yellow, cyan, and magenta for (a) and yellow for (b and c)),  
712 Africa (blue), America (red), and Europe (star). (a) A core SNP-based maximum  
713 likelihood phylogeny of 469 genomes with geographical origins highlighted. The tree  
714 was rooted on *B. pseudomallei* MSHR5619, the most genetically distant isolate based  
715 on pairwise SNP distance (see methods and Supplementary Figure 10). The outer ring  
716 represents population clusters based on BAPS hierarchical clustering (Group 1 – 19).  
717 Apart from Group 15, which is paraphyletic and marked by two black arrows, other  
718 groups each form a monophyletic branch. (b) Pan-genome accumulation curve  
719 representing rates of new gene discovery in isolates collected from different  
720 geographical origins. The order of new genome added was permuted 1,000 times to  
721 accommodate possible assortment. (c) Summary of core and accessory genomes of  
722 isolates grouped by geographical origins.

723

### 724 **Figure 2 Timeline of trans-continental and sub-regional spread of *B.*** 725 ***pseudomallei***

726 (a) Estimated time when the most recent common ancestor (MRCA) of each cluster  
727 emerged. Time (black dots) and 95% highest posterior density (horizontal line) were  
728 estimated by BEAST for those clusters with temporal signals. Estimations were  
729 performed separately for chromosome I (solid lines), and II (dotted lines).  
730 Overlapping estimations between the two chromosomes provide further confidence in  
731 the time interval in which the MRCA emerged. The estimation for chromosome I of  
732 group 4 did not reach a credible effective sample size and was excluded. (b)  
733 Transatlantic slave trade routes and sampling locations of African and American  
734 isolates. Each dot represents the geographical origin of isolates used for the time  
735 estimation with the size proportional to the number of isolates. (c) The geographical  
736 landscape and isolates used to determine sub-regional connectivity. Isolates  
737 representing six Southeast Asian countries were plotted on the map, highlighting the  
738 geographical proximity of the Mekong group, and the Malay group. The number of  
739 isolates from each country was annotated.

740

### 741 **Figure 3 Region-specific genetic signatures**

742 Functional categories of genes (COG) localised in region-specific loci. One-sided  
743 Fisher's exact test was used to search for terms that showed significant departure from  
744 random expectation in the reference genome. Asterisks highlight terms with  
745 heightened frequency following Bonferroni correction for multiple testing. \*denotes  
746 terms with p-value  $<10^{-9}$ , while \*\* denotes terms with p-value  $<2.2 \times 10^{-16}$ .