

1 **Genetic diversity, determinants, and dissemination of *Burkholderia pseudomallei***
2 **lineages implicated in melioidosis in northeast Thailand**

3 Rathanin Seng¹, Chalita Chomkatekaew^{2,3}, Sarunporn Tandhavanant¹, Natnaree Saiprom¹,
4 Rungnapa Phunpang¹, Janjira Thaipadungpanit⁴, Elizabeth M Batty^{2,5}, Nicholas PJ Day^{2,5},
5 Wasun Chantratita⁶, T. Eoin West^{1,7,8}, Nicholas R Thomson⁹, Julian Parkhill³, Claire
6 Chewapreecha^{2,4,9,10*}, Narisara Chantratita^{1,2*}

7 ¹ Department of Microbiology and Immunology, Faculty of Tropical Medicine, Mahidol University,
8 Bangkok, Thailand.

9 ² Mahidol Oxford Tropical Medicine Research Unit, Faculty of Tropical Medicine, Mahidol
10 University, Bangkok, Thailand.

11 ³Department of Veterinary Medicine, University of Cambridge, UK

12 ⁴Department of Clinical Tropical Medicine, Faculty of Tropical Medicine, Mahidol University,
13 Bangkok, Thailand.

14 ⁵Centre for Tropical Medicine and Global Health, Nuffield Department of Medicine, University of
15 Oxford, Oxford, UK.

16 ⁶Center for Medical Genomics, Faculty of Medicine Ramathibodi Hospital, Mahidol University,
17 Bangkok, Thailand

18 ⁷Division of Pulmonary, Critical Care & Sleep Medicine, Department of Medicine, University of
19 Washington, Seattle, Washington, USA.

20 ⁸ Department of Global Health, University of Washington, Seattle, Washington, USA.

21 ⁹Parasites and Microbes, Wellcome Sanger Institute, Cambridge, UK.

22 ¹⁰ Previous Affiliations: Bioinformatics and Systems Biology Program, School of Bioresource
23 and Technology, King Mongkut University of Technology Thonburi, Bangkok, Thailand.

24

25 *Contribute equally

26 For sequencing, and data analysis - correspond to Claire Chewapreecha

27 (claire@tropmedres.ac)

28 For sample collection, sampling processing and metadata - correspond to Narisara Chantratita

29 (narisara@tropmedres.ac)

30

31 **Abstract**

32 Melioidosis is an often-fatal neglected tropical disease caused by an environmental bacterium
33 *Burkholderia pseudomallei*. However, our understanding of the disease-causing bacterial
34 lineages, their dissemination, and adaptive mechanisms remains limited. To address this, we
35 conducted a comprehensive genomic analysis of 1,391 *B. pseudomallei* isolates collected from
36 nine hospitals in northeast Thailand between 2015 and 2018, and contemporaneous isolates
37 from neighbouring countries, representing the most densely sampled collection to date. Our
38 study identified three dominant lineages with unique gene sets enhancing bacterial fitness,
39 indicating lineage-specific adaptation strategies. Crucially, recombination was found to drive
40 lineage-specific gene flow. Transcriptome analyses of representative clinical isolates from each
41 dominant lineage revealed heightened expression of lineage-specific genes in environmental
42 versus infection conditions, notably under nutrient depletion, highlighting environmental
43 persistence as a key factor in the success of dominant lineages. The study also revealed the
44 role of environmental factors – slope of terrain, altitude, direction of rivers, and the northeast
45 monsoons - in shaping *B. pseudomallei* geographical dispersal. Collectively, our findings
46 highlight persistence in the environment as a pivotal element facilitating *B. pseudomallei* spread,
47 and as a prelude to exposure and infection, thereby providing useful insights for informing
48 melioidosis prevention and control strategies.

49

50 **Introduction**

51 Melioidosis, a severe infectious disease, affects an estimated 165,000 cases globally each year,
52 of which 89,000 are fatal¹. The disease is caused by *Burkholderia pseudomallei*, a Gram-
53 negative bacillus found in soil and contaminated water across tropical and sub-tropical regions.
54 Historically, limited access to microbiology laboratories for culture-confirmed diagnosis led to
55 underreporting, particularly in lower- and middle-income countries². However, improved
56 infrastructure and awareness have led to increases in reported cases across South Asia,
57 Southeast Asia, East Asia³⁻⁸ and Australia^{9,10}. In Southeast Asia, the disease incidence is often
58 linked to agriculture practice, particularly during the rainy seasons when rice paddy fields are
59 flooded for planting. The flooded terrain enables the bacterium in the soil to surface, potentially
60 exposing farmers to *B. pseudomallei* and subsequently leading to melioidosis¹¹. Additionally,
61 many cases of melioidosis have been associated with severe weather events¹²⁻¹⁴. While climate
62 likely influences human encounters with *B. pseudomallei*, further investigation is needed to fully
63 understand the mechanisms linking environmental factors to melioidosis epidemiology.

64

65 Understanding the population structure, dissemination and adaptation of *B. pseudomallei* in
66 these climatically-challenged endemic regions requires a large-scale, geographically and
67 chronologically densely-sampled, genetic dataset. Previous studies, albeit limited in sample
68 size, have demonstrated that *B. pseudomallei* dissemination is driven by both anthropogenic
69 and environmental factors^{15–18}. Streams^{19,20}, monsoons, typhoons and cyclones^{13,14,21,22} were
70 identified as significant contributors to bacterial dissemination, highlighting the importance of
71 bacterial persistence across a range of environmental conditions. *B. pseudomallei* exhibits
72 remarkable survival capabilities across diverse environments, spanning from wet to dry,
73 nutrient-depleted soil^{23–28} thereby enabling the bacterium to thrive in various ecological niches.
74 Previous studies have noted the temporal and geographical co-existence of multiple *B.*
75 *pseudomallei* lineages^{15,16,29}. However, little is known about their distinct genetic content and
76 adaptive strategies. Identification of lineage-specific genes associated with bacterial persistence
77 and disease escalation will be essential to develop disease control strategies.

78
79 In this study, we conducted a population genomics analysis using combined *B. pseudomallei*
80 isolates from melioidosis patients across nine provinces in the northeast Thailand including
81 Buriram, Khon Kaen, Mahasarakam, Mukdahan, Nakhon Phanom, Roi-Et, Sisaket, Surin and
82 Udon Thani⁸; totaling 1,265 isolates collected from July 2015 to December 2018. Additionally,
83 we incorporated contemporary environmental and clinical collections from Thailand and
84 neighbouring countries^{15,29–34}, consisting of 15 clinical isolates and 111 environmental isolates
85 (**Figure 1a, Supplementary data 1**). Our comprehensive analysis, including a total of 1,391
86 isolates, revealed the population structure, dissemination patterns, and genetic diversity of this
87 bacterium. We identified genetic determinants associated with dominant lineages and
88 investigated their biological functions and expression conditions in three representative isolates,
89 each representing a dominant lineage. This provides insights into the strategies employed by *B.*
90 *pseudomallei* lineages for successful persistence in the environment, ultimately leading to
91 human exposure and infection.

92

93 **Results**

94 **Population structure analysis revealed the successful *B. pseudomallei* lineages and** 95 **mixture of clinical and environmental isolates**

96 To define the population structure of clinical and environmental *B. pseudomallei* in a
97 hyperendemic area of northeast Thailand and neighbouring regions (n = 1,391), we performed
98 four independent approaches. PopPUNK³⁵ analysis was performed on genome assemblies

99 **(Supplementary data 2)**. Additionally, we constructed three maximum-likelihood (ML)
100 phylogenies³⁶, each based on different sets of single nucleotide polymorphisms (SNPs): core
101 genomes (n= 77,156 SNPs), core gene multilocus typing³⁷ (cgMLST, n = 46,945 SNPs), and
102 seven-gene multilocus typing genes³⁸ (MLST, n = 31 SNPs). These approaches facilitated the
103 grouping of isolates with close genetic similarity into distinct lineages. Notably, both PopPUNK
104 analysis and core genome SNPs phylogeny yielded consistent results **(Supplementary Figure**
105 **1)** clustering the population into three dominant lineages **(Figure 1b)**. The average pairwise
106 core SNP distance within each dominant lineage was 549, 351, and 517 SNPs for lineage 1, 2,
107 and 3, respectively, in contrast to the pairwise core SNP distance of 1,087 SNPs within the total
108 population **(Figure 1c)**. This lower pairwise core SNP distance across lineages confirmed the
109 genetic relatedness as defined by PopPUNK and core genome SNP phylogenetic analysis.
110 While cgMLST displayed conservation for two out of three dominant lineages, MLST exhibited
111 inconsistencies across dominant lineages with lower phylogenetic resolution and poorer
112 bootstrap support compared to other methodologies **(Supplementary Figure 1)**. Consequently,
113 we relied on the population delineated by PopPUNK and core genome SNP phylogeny for
114 subsequent investigations.

115
116 The three predominant lineages (denoted as lineage 1 to 3) comprised 312, 297, 125 isolates,
117 respectively. They accounted for 52.8% of the studied population and persisted throughout the
118 sampling period. Interestingly, each lineage peaked during the rainy season, correlating
119 agricultural practices at the onset of rainfall with increased environmental exposure and
120 subsequent melioidosis infections **(Figure 1d)**. Despite the small sample size of environmental
121 isolates, we observed a clustering of these isolates with clinical isolates within each dominant
122 lineage, indicating their core genetic similarities and shared origin. The ratio of environmental to
123 clinical isolates varied across lineages (Chi-square test with Monte Carlo resampling p-value
124 5.00×10^{-4} , **Supplementary Figure 2**). Due to the substantially lower number of environmental
125 isolates used and incomplete geographical distribution matching between clinical and
126 environmental isolates, caution is warranted in interpreting these results. Nevertheless, our
127 findings highlight a mixing of environmental and clinical samples, suggesting that clinical
128 isolates could serve as a surrogate for tracking the dissemination of an environmental
129 bacterium, especially in the absence of equally comprehensive environmental samples.

130
131 **Genetic evidence identifies patterns of *B. pseudomallei* dissemination in Northeast**
132 **Thailand**

133 We examined the dissemination patterns of *B. pseudomallei* in northeast Thailand. Except for
134 Maharakam where samples were limited, all three dominant lineages were present across the
135 rest of eight studied provinces (**Figure 1a**). This prompted us to focus the analysis on these
136 dominant lineages to identify consistent geographical distributions underlying their spread in the
137 region. We generated lineage-specific phylogenies to improve genetic resolution for
138 transmission analysis. Additionally, we reconstructed ancestral histories of provincial origins,
139 quantified the number of inter-provincial transmissions to examine transmission patterns, and
140 estimated the time of the most recent common ancestor of each dominant lineage and its sub-
141 lineages (**Supplementary Figures 3 – 5**, see Methods). This allowed us to link the emergence
142 of these lineages to historical events that might have affected their transmission dynamics.
143 While transmission signals could reflect distinct dissemination patterns in each dominant lineage
144 or its sub-lineages, we also considered the possibility of shared factors leading to a uniform
145 geographical distribution. Notably, we observed consistent dissemination patterns in 14 out of
146 28 provincial pairs across the three dominant lineages (**Figure 2a, Supplementary Figure 6**).
147 Eight out of these 14 provincial pairs potentially correlated with the slope of terrain altitude
148 between provinces or the natural flow of rivers in the region (**Figure 2b**). These pairs included
149 “Udon Thani-to-Khon Kaen”, “Udon Thani-to-Buriram”, “Udon Thani-to-Mukdahan”, “Udon
150 Thani-to-Surin”, “Khon Kaen-to-Buriram”, “Nakhon Phanom-to-Mukdahan”, “Roi-Et-to-Surin” and
151 “Surin-to-Sisaket”. Northeast Thailand is described as a saucer-shaped plateau, with elevations
152 ranging from over 200 meters above sea level in the northwestern corner (parts of Udon-Thani,
153 and Khon-Kaen) to less than 100 meters in the southeast (parts of Mukdahan and Sisaket), and
154 gradually descending toward the Mekong River in the east^{39,40}. The Mun River originates from
155 elevated hills in central Thailand, streaming eastward through Buriram, Surin, Sisaket before
156 merging with the Mekong River. Similarly, the Chi River, a tributary to the Mun, also originates
157 from the central Thailand mountains. The Chi flows eastward through Khon Kaen,
158 Maharakam, Roi-Et and converges with the Mun in Sisaket⁴⁰.

159 Additionally, another set of three out of 14 conserved patterns coincided with the wind direction
160 of the northeastern monsoon during the dry season (“Nakhon Phanom-to-Buriram”, “Mukdahan-
161 to-Buriram”, and “Mukdahan-to-Surin”). Thailand experiences two predominant monsoon
162 seasons: the southwest monsoon from May to October and the northeast monsoon from
163 November to April. The southwest monsoon brings heavy rainfall from the southwest to
164 northeast, often marking the start of the agriculture season (**Figure 1c**)^{41,42}. Conversely, the
165 northeast monsoon brings dry winds from the northeast to southwest. Our observation implies

166 the potential for dry winds to transport aerosolised soil contaminated with *B. pseudomallei*
167 westward. Despite that previous air sampling during Thailand's rainy season did not detect *B.*
168 *pseudomallei*, exploring the impact of dry winds during the northeastern monsoon is essential.
169 This consideration becomes even more pertinent given reports of the long-range transport of
170 particles and small organic matters, such as PM2.5 and PM10, via the northeast monsoon
171 elsewhere in Southeast Asia^{43,44}. Furthermore, we successfully estimated the time of most
172 recent of ancestry for a sub-lineage 1.3, a descendant of lineage 1. Our finding revealed that
173 this sub-lineage emerged around 2011 (95% HPD of 2000-2014, **Supplementary Figure 3**).
174 The age of this sub-lineage implies that older lineages, such as its parent lineage 1 likely
175 experienced multiple monsoon seasons, which possibly resulted in the observed patterns. Apart
176 from the terrain slope, inland rivers, canal systems⁴⁵ and regular monsoons^{41,42}; various factors
177 likely contributed to shaping *B. pseudomallei* dissemination in northeast Thailand.
178 Anthropogenic activities, such as human migration between the provinces, may also contribute
179 to the observed pattern; however without access to comprehensive human movement data, this
180 aspect remains challenging to investigate.

181 **Genetic markers potentially contributing to the emergence of successful *B. pseudomallei*** 182 **lineages**

183 The co-existence of multiple *B. pseudomallei* lineages within the same geographical areas and
184 timeframe implies the presence of diverse adaptive strategies which enable them to thrive in a
185 shared ecological niche. While some smaller lineages may be sporadically detected, the
186 persistence of the three dominant lineages throughout the sampling period supports their fitness
187 and successful adaptive strategies in this niche. We next sought to identify genes that were
188 present in isolates that form each dominant lineage, or its sub-lineage; but absent in non-
189 dominant lineages (see Methods). Out of total 15,237 genes in the pan-genome outlined from
190 this population (see Methods), 5,577 genes were conserved across the entire population while
191 9,660 genes were variably present (accessory genes). Dominant lineage-specific genes were
192 defined as accessory genes present in $\geq 95\%$ of isolates within any of the dominant lineages or
193 their sub-lineages, but present in $\leq 15\%$ of isolates outside these lineages. Among these, 247
194 genes were identified as lineage-specific with their specificity to each dominant lineage and sub-
195 lineage tabulated in **Supplementary data 3**. The majority of dominant lineage-specific genes
196 were poorly characterised and annotated as hypothetical proteins (**Figure 3**). To gain insights
197 into the potential functions, we annotated them using Gene Ontology (GO terms)⁴⁶ which
198 classify them by Biological process, Molecular function, and Cellular component (Figure 3b, see

199 Methods). Of the 247 dominant lineage genes, GO terms could be assigned to 27 genes for
200 Biological Process, 68 for Molecular Function, and 12 for Cellular component. For genes that
201 could be assigned GO terms, functions involved in “DNA integration”, “DNA recombination”, and
202 “DNA methylation” might indicate their potential roles in horizontal gene acquisition and
203 protection against incoming foreign DNA through site-specific DNA methylation. Furthermore,
204 GO terms associated with “DNA binding” and “Regulation of DNA-templated transcription” may
205 suggest lineage-specific regulation of the expression of these genes.

206

207 **Lineage-specific genes were selectively expressed**

208 To delve deeper into the functionality of dominant lineage-specific genes, we explored their
209 expression patterns across both environmental and infection conditions. We selected
210 representative strains - “K96243” (lineage 1 – sub-lineage 1.1), “UKMD286” (lineage 2 – sub-
211 lineage 2.1), and “UKMH10” (lineage 3 – sub-lineage 3.2) - all are clinical isolates with pre-
212 existing gene expression profiles under infection and environmental conditions⁴⁷⁻⁴⁹. For
213 environmental conditions, K96243 was exposed to water⁴⁷, while UKMD286 and UKMH10 were
214 cultivated in a soil extract medium^{48,49} to mimic *B. pseudomallei* in the environment. For
215 infection conditions, K96243 and UKMD286 were used in murine challenges^{47,48} and isolated
216 from mice organs, while UKMH10 was subjected to human plasma⁴⁹ to simulate host infection.
217 This approach facilitated the comparison of differentially expressed lineage-specific genes
218 between environmental and infection conditions. Although each representative strain carried a
219 complete set of lineage-specific genes for their respective sub-lineages (K96243 with 47 genes,
220 UKMD286 with 27 genes, and UKMH10 with 14 genes), their collective representation
221 accounted for 69 out of 247 total lineage-specific genes (27.9%) due to observed genetic
222 diversity within the dominant lineage. Notably, 11 out of 47 lineage-specific genes in K96243
223 and 6 out of 27 lineage-specific genes in UKMD286 were up-regulated in the environmental
224 conditions (**Figure 4a to 4c, Supplementary data 4**). Notably, none of the lineage-specific
225 genes showed up-regulation during the infection condition. The remaining lineage-specific
226 genes did not exhibit preferential expression in either environmental or infection conditions. The
227 elevated expression level of lineage-specific genes in the environmental condition was
228 unexpected considering that all representative strains were clinical isolates. This observation
229 potentially suggests that dominant lineage-specific genes may play more substantial roles in
230 bacterial environmental survival than in host pathogenicity.

231

232 To thrive in its environmental habitat, *B. pseudomallei* must cope with ranges of physical,
233 chemical and biological stresses such as desiccation, temperature fluctuations, osmotic
234 changes, oxidative stress, UV exposure, nutrient scarcity, changes in pH, exposure to heavy
235 metals, competition from antibiotics released by other microbes, and predation by eukaryotes.
236 We leveraged the extensive condition-wide transcriptome data spanning 62 conditions available
237 for K96243⁴⁷, a representative strain from lineage 1, to compare the expression patterns of
238 lineage-1-specific genes with the rest of genes in the K96243 genome (**Figure 4d**). Our analysis
239 revealed that lineage-1-specific genes within K96243 exhibited a higher level of gene
240 expression when the bacterial cell experienced nutrient deprivation compared to other genes in
241 K96243 genome (Two-sided Fisher's exact test p-value = 9.23×10^{-5}). This finding implies that
242 lineage 1 might possess an adaptive strategy to persist in nutrient-depleted soil, which is not
243 uncommon in melioidosis endemic areas^{27,28,50}, before being acquired by a human host and
244 subsequently causing the disease.

245

246 **Example of lineage-specific genes**

247 The majority of lineage-specific genes were located within genomic islands (GI)^{30,51,52}. These
248 regions are characterised by anomalies in %G+C content or dinucleotide frequency signatures,
249 or the presence of genes associated with mobile genetic elements such as insertion sequence
250 (IS) elements and bacteriophages. Notably, we observed that a cluster of genes specific to
251 lineage 1 (*BPSS2060* to *BPSS2072*) formed a mosaic structure within a putative metabolic
252 island known as GI 16³⁰. Although several variations of GI 16 have been reported (GI16, GI16.1,
253 GI16.2, GI16a, GI16b, and GI16b.1)⁵¹, it typically spans 60 kb (*BPSS2051* to *BPSS2090*) and
254 carries several known virulence determinants and genes that enhance metabolic versatility.
255 While certain virulence factors, such as the filamentous haemagglutinin (*BPSS2053*) required
256 for host cell adhesion and its processing protein were conserved across multiple lineages
257 observed in our study, genes encoding functions that potentially expand the metabolic repertoire
258 were specific to dominant lineages. For example, the mosaic structure of GI 16 (*BPSS2060* to
259 *BPSS2072*), specific to lineage 1 (**Supplementary data 3**), contains genes involved in
260 alternative nutrient catabolism and anabolism (*BPSS2060*, *BPSS2065*, *BPSS2067*, *BPSS2068*,
261 and *BPSS2072*), transcriptional regulation (*BPSS2061*), and substrate transport (*BPSS2064*,
262 *BPSS2071*). Out of 11 lineage-specific genes located in the mosaic structure of GI 16, eight
263 were found to be upregulated during the early phase of nutrient starvation while remaining silent
264 during infections. This finding reflects the functional division of GI16 where its lineage-specific
265 mosaic structure contains genes that contribute to metabolic versatility, while its core structure

266 encodes virulent determinants associated with disease implications. It is important to note that
267 the structure of GI 16 may vary across different regions due to the plasticity of genomic islands
268 and changes in selection pressures. While this observation is significant for a dominant lineage
269 in northeast Thailand, it may not be generalisable to other geographical locations.

270

271 **Lineage-specific genes were introduced by homologous recombination**

272 Homologous recombination has been shown to play a significant role in facilitating the
273 acquisition and loss of genes, and the generation of mosaic structures within the GI of *B.*
274 *pseudomallei*⁵³. To better understand its association with lineage-specific genes, we identified
275 recombination events and quantified the rates of recombination in the dominant lineages. The
276 ratio of polymorphisms introduced through recombination compared to those introduced by
277 mutation (r/m) was 3.7, 4.6 and 2.2 for lineages 1, 2, and 3 respectively (**Table 1**). A very high
278 proportion of genes underwent recombination at least once: 99.5% of genes in lineage 1, 99.9%
279 in lineage 2, and 96.6% in lineage 3. Furthermore, every lineage-specific gene within each
280 dominant lineage underwent recombination (**Supplementary Figure 8**). The bacterial restriction
281 modification (RM) systems prevent the invasion of foreign DNA and restrict gene flow between
282 *B. pseudomallei* lineages³¹. Notably, components of this system including a type I restriction
283 system and modification methylase were among dominant lineage-specific genes (*BPSL0947-*
284 *BPSL0948* in lineage 1, and their homologues in lineage 2 and 3). They may act as a barrier for
285 homologous recombination and potentially modulate lineage-specific genetic diversity. This
286 highlights the intricate interplay between recombination, lineage-specific genes, and the RM
287 system in shaping the genetic landscape of *B. pseudomallei* in northeast Thailand.

288

289 **Discussion**

290 Our analysis of *B. pseudomallei* population genomics enhances our understanding of the
291 evolution and adaptive strategies employed by the dominant lineages in the melioidosis
292 hyperendemic region of northeast Thailand and neighbouring countries. Through an
293 unprecedentedly dense sampling effort between 2015 and 2018, we were able to determine the
294 co-existence of three dominant lineages, characterise their dissemination patterns, and identify
295 lineage-specific genes possibly contributing to their success during the studied period. By
296 analysing transcriptome data from representative strains of each dominant lineage, we gained
297 further insights into their adaptive strategies, particularly emphasising bacterial persistence in
298 the environment as crucial for subsequent host acquisition and infection. Nevertheless, our
299 study has a few limitations.

300

301 Due to limited environmental surveillance and the scarcity of environmental isolates in
302 Southeast Asia, our study primarily relied on clinical isolates. Nonetheless, the co-occurrence of
303 clinical and environmental isolates within the same lineage, coupled with the direct acquisition of
304 clinical isolates from the environment, suggests that our findings may hold broader implications
305 for environmental isolates. Furthermore, the lineage classification report in our study is subject
306 to potential alterations over time with the introduction of new data. As *B. pseudomallei*
307 continuously adapts to environmental pressures, the composition of lineages 1, 2, and 3, along
308 with their respective lineage-specific genes may shift. New lineages with superior fitness or
309 selective advantages, of different adaptive strategies could potentially outcompete the existing
310 dominant lineages. Consequently, this may lead to emergence of new lineage classifications in
311 the future. Continued surveillance efforts including both clinical and environmental samples will
312 be essential. This ongoing monitoring will be pivotal in identifying alterations within the bacterial
313 population, uncovering new adaptive strategies, and evaluating their impact on disease
314 dynamics over time.

315

316 Our understanding of the significance of bacterial persistence as a strategy for successful
317 lineages stems from the observed up-regulation of dominant lineage-specific genes in
318 environmental conditions, along with the increased expression of lineage-1-specific genes
319 during nutrient deprivation. However, it is essential to recognise the genetic diversity within each
320 dominant lineage. The representative strains used in our study carried a subset of the lineage-
321 specific genes corresponding to their lineage. As a result, lineage-specific genes absent in
322 these strains, often annotated as hypothetical proteins, remained unexplored in our analysis.
323 Moreover, our detailed characterisation of lineage-1-specific genes using a comprehensive
324 transcriptome dataset of 62 distinct conditions might not capture the full spectrum of conditions
325 encountered by *B. pseudomallei* in its natural habitat. There is a possibility that other adaptive
326 strategies were overlooked in this study, indicating scope for future exploration. Despite these
327 shortcomings, our dataset and analysis currently represent one of the most comprehensive
328 efforts to date. Future research will prioritise the generation of a more extensive condition-wide
329 transcriptome, covering a broader range of conditions and incorporating strains with diverse
330 genetic variations to identify other adaptive strategies employed by *B. pseudomallei*.

331

332 Our findings underscore the significance of environmental persistence in driving the success of
333 dominant lineages 1 and 2, notably highlighting lineage-1-specific genes in mediating bacterial

334 survival under nutrient depletion. It remains uncertain whether lineage 3 adopts a similar
335 strategy. Nevertheless, our results align with previous soil sampling studies, which consistently
336 observed a higher prevalence of *B. pseudomallei* in nutrient-depleted compared to nutrient-rich
337 soil^{28,50}. Additionally, a molecular evolutionary study also supports the species' long-term
338 adaptation to survive nutrient scarcity²⁷. The northeast region of Thailand, where our samples
339 were primarily collected, has inherently low-fertility soil. This is exacerbated by intensified
340 agriculture, monoculture, excessive synthetic fertilizer use, and poor land management,
341 resulting in depleted soil nutrients and organic matter⁵⁴. This presents a challenging
342 environment for *B. pseudomallei* to thrive, thereby potentially selecting for successful lineages
343 with persistent traits as observed in our study.

344
345 Our analyses also highlight the role of various factors such as differences in terrain altitude^{55,56},
346 river flow dynamics⁴⁰, and the northeast monsoon⁴¹ contribute to shaping the dissemination
347 patterns of *B. pseudomallei*. These drivers of dissemination are influenced by both natural and
348 human activities. For instance, strong winds can carry dried soil particles⁵⁷, potentially
349 containing *B. pseudomallei* over distances. Climate change-induced alterations in vegetation
350 cover might expose soil to rainfall and winds⁵⁸, impacting the bacterial spread. Additionally,
351 deforestation can disrupt natural barriers like trees and shrubs, accelerating water runoff^{57,58} and
352 potentially facilitating the wide-ranging dissemination of *B. pseudomallei* during flood.
353 Considering these dynamics, the strategy of bacterial persistence likely plays a pivotal role in its
354 widespread dissemination within the region, thereby influencing disease prevalence. Therefore,
355 an effective disease control strategy should integrate both environmental and clinical public
356 health measures to effectively mitigate the impact of melioidosis.

357 **Materials and Methods**

358 **Data collection and bacterial isolates**

359 The *B. pseudomallei* isolates in our study included a cohort⁸ from northeast Thailand gathered
360 between July 2015 to December 2018 consisting of 1,265 clinical isolates. We also incorporated
361 a contemporaneous dataset from Thailand and neighbouring regions^{15,29–34}, comprising 15
362 clinical and 111 environmental isolates sourced from previous publications. In total, 1,391 *B.*
363 *pseudomallei* genomes were used in this study. Their metadata and accession numbers were
364 documented in **Supplementary data 1**.

365
366 The northeast Thailand collection was collected from patients participated in our longitudinal
367 cohort study⁷. The patients were from nine provinces including Udon Thani (n = 230), Mukdahan
368 (n = 198), Roi Et (n = 195), Surin (n = 170), Nakhon Phanom (n = 135), Buriram (n = 123),
369 Sisaket (n = 107), Khon Kaen (n = 96) and Maha Sarakham (n = 11), who were admitted to nine
370 hospitals included in our cohort. The ethical approval for the cohort study was obtained from the
371 Ethics Committee of the Faculty of Tropical Medicine, Mahidol University (MUTM 2015-002-001
372 and MUTM 2021-055-01). The isolates were obtained from various clinical samples, including
373 blood (71.8%), pus (12.5%), sputum (9.9%), body fluid (3.5%), urine (1.9%) and tissue (0.4%)
374 (**Supplementary data 1**). As latent infection accounts for < 5% of the cases, the majority of
375 clinical cases likely directly acquired from the environment⁸. The numbers of enrolled cases
376 were lower at the beginning of the study due to the delayed sample collection in some study
377 sites, resulting in inconsistent number of bacterial isolates used across different sites. When
378 applicable, a permutation test was performed to ensure that an unequal number of isolates did
379 not impact the temporal or spatial analysis.

380 381 **Culture confirmation of *B. pseudomallei*, DNA extraction and whole genome sequencing**

382 All 1,265 *B. pseudomallei* samples from the northeast Thailand collection were cultured on
383 Ashdown, selective agar plates and confirmed the species using latex agglutination test and
384 matrix-laser absorption ionisation mass spectrometry (MALDI-TOF MS). A single colony from
385 Ashdown agar plate was subjected to culture in Luria-Bertani (LB) broth and subsequently used
386 for DNA extraction. Genomic DNA was extracted using QIAamp DNA Mini Kit (Qiagen,
387 Germany). All genomic DNA were processed for the 150-base-read library preparation and
388 sequencing using Illumina HiSeq2000 system with 100-cycle paired-end runs at Wellcome
389 Sanger Institute, Cambridge UK. An average of 71X read depth was achieved. To control the
390 potential contamination in each sample with other closely related species, we assigned

391 taxonomic identity using Kraken⁵⁹ v.1.1.1. We then estimated the genome completeness and
392 species confirmation using CheckM⁶⁰ v.1.2.2 and FastANI⁶¹ v.1.31, respectively. The quality
393 control data of the 1,265 genomes were listed in **Supplementary data 2**.

394

395 **Genome assembly and mapping from short read data**

396 Short reads were *de novo* assembled using Velvet v.1.2.10⁶² followed by optimisation as
397 previously described²⁹, with their quality scores reported in **Supplementary data 2**. Short reads
398 were also mapped against several reference genomes, including a strain K96243³⁰ (accession
399 numbers BX571965 and BX571966) to determine the whole population structure, and lineage-
400 specific references to improve the resolution for lineage-specific analyses. We selected
401 K96243³⁰ genome as a population-wide reference due to its origin in northeast Thailand,
402 aligning with the geographical focus of our study. Additionally, its well-characterised and
403 complete genome further streamlined subsequent analyses. For all mapping, variants were
404 called using Snippy v.4.6.0 (<https://github.com/tseemann/snippy>). To avoid mapping errors and
405 false SNPs, we filtered out SNPs covered by less than 10 reads and found in a frequency of
406 less than 0.9.

407

408 **Defining whole population structure**

409 **PopPUNK clustering**

410 PopPUNK v.2.6.0³⁵ was run on 1,391 assembled genomes. To define the core and accessory
411 distance between each pair of isolates, the assemblies were hashed at different k-mers. The
412 population model was fit using command line “poppunk-runner.py --fit-model --distance
413 <database.dists> --output <database> --full-db --ref-db <database> --min-kmer 15 --max-kmer
414 31 --max-a-dist 0.53 --K 4 --k-step 2” with the result density of 0.028, transitivity of 0.992, and
415 network score of 0.8961.

416

417 **Maximum likelihood phylogenies from core genome SNP, cgMLST and MLST**

418 An alignment of full genome was created by mapping whole genome sequences of each *B.*
419 *pseudomallei* against a complete genome of K96243³⁰ strain. From this alignment, 4,221
420 cgMLST loci based on a scheme described in³⁷ were extracted and concatenated to form
421 cgMLST alignment. Additionally, seven MLST loci, as per scheme described in³⁸ were
422 extracted from the same alignment and concatenated to create the MLST alignment. Core
423 genome SNP alignment was identified from a full genome alignment using snp-sites⁶³ v.2.5.1,
424 with genomic islands⁵¹ masked. Separate maximum likelihood phylogenies were constructed for

425 core genome SNP alignment, cgMLST alignment, and MLST alignment using IQ-TREE³⁶
426 v.2.0.3. Standard model selection in IQ-TREE determined the best-fit model as
427 TVM+F+ASC+R6 for all three phylogenies. To assess the robustness of the phylogenetic trees,
428 a 1,000 bootstrap support was performed for each tree.

429

430 **Comparison of population structure outlined phylogeny constructed from core genome** 431 **SNPs, cgMLST, MLST and PopPUNK**

432 To test for consistency between phylogenetic trees constructed from core genome SNPs,
433 cgMLST and MLST alignment, we use the R package treespace⁶⁴ v. 1.1.4.3 to explore the tree
434 tip distributions. We compared pairwise tree distances within the first 100 bootstraps within each
435 alignment category (indicative of bootstrap support strength), and across trees generated from
436 different alignment categories (indicating proximity between tree categories). The tree pairwise
437 distances were computed, and principal components (PCs) were derived with eigenvalues
438 calculated for different PCs. The similarity among phylogenies from each alignment category
439 was assessed using two PCs dimensions, which jointly accounted for >90% of variability in
440 pairwise distance (**Supplementary Figure 1a**). The scatter plot of PCs revealed a close
441 clustering of bootstrap trees from core genome SNP and cgMLST, while the bootstrap trees
442 from MLST alignment showed greater dispersion, highlighting less consistency in the trees
443 generated by the MLST approach. We further compared the consistency between the median
444 phylogenetic tree of each alignment category and PopPUNK classification was visually
445 compared using iTOL⁶⁵ (**Supplementary Figure 1b**).

446

447 **Specific lineage analysis**

448 To investigate the dissemination and genetic diversity within each dominant lineage, we
449 conducted individual genome alignments, recombination removal, and maximum likelihood
450 phylogeny. To enhance the sensitivity of variant, we selected closely related genomes as
451 references for each lineage. Specifically, the complete genome *B. pseudomallei* strain K96243
452 (accession numbers BX571965 and BX571966) served as the mapping reference for lineage 1,
453 while new reference genomes were created for lineages 2 and 3.

454

455 For lineage 2, we chose a representative isolate 27035_8#57 and subjected it to long-read
456 sequencing on a local MinION sequencer following the manufacturer's standard protocol
457 (Oxford Nanopore Technologies, Oxford, United Kingdom). A complete hybrid assembly of the
458 long-read and short-read sequence data of this strain was performed using Unicycler⁶⁶ v.0.8.4.

459 The resulting hybrid assembly of the 27035_8#57 genome was employed as the mapping
460 reference for this lineage.

461
462 In the absence of representative complete genomes for lineage 3, we selected the best quality
463 de novo assembly of isolate 27035_8#119 and orientated its contigs according to strain K96243
464 using ABACAS⁶⁷ v.1.3.1. This genome was used as a mapping reference for lineage 3.

465
466 For lineage-specific mapping, Snippy v.4.6.0 as employed as in the whole population analysis.
467 All genome alignments were subjected to Gubbins⁶⁸ v.3.1.3, a recombination identification tool,
468 to detect and remove recombination fragments. This process determined the genetic diversity
469 introduced by horizontally acquired elements and vertically inherited SNPs, thereby producing
470 recombination-free SNP alignments for phylogenetic reconstruction. Maximum-likelihood
471 phylogenies were constructed using recombination-free SNP alignment of each dominant
472 lineage using IQ-TREE³⁶ v.2.0.3 with TVM+F+ASC+R6 and 1,000 replicates of bootstrap
473 support. The overall proportion of nodes with $\geq 80\%$ bootstrap support of lineage-specific
474 phylogenies reached 83.5%.

475 476 **Dating the timeline for lineages and sub-lineages**

477 To enable dating analysis, we further divided each lineage into sub-lineages using R package
478 rhierbaps⁶⁹ v.1.1.4. We inferred evolutionary timeline and estimated the age of each lineage and
479 sub-lineage based on the isolate's collection date. A Bayesian molecular dating provided in R
480 package BactDating⁷⁰ v.1.1.1. was employed to assess the temporal signals by examining a
481 positive correlation between the isolate's collection date and the root-to-tip distance.
482 Recombination removed phylogenies were used in this analysis. A date-randomisation test,
483 consisting of 100 permutations, was performed to assess the robustness of the temporal signal
484 compared to noise.

485
486 Notably, the temporal signals were discernable at the sub-lineage level rather than the broader
487 lineage level. Among the 10 sub-lineages, only one sub-lineage (lineage 1.3) exhibited a
488 positive correlation in their clock signals. Given the limited sample size of lineage 1.3, we
489 employed a strict clock model to prevent parameter over-fitting. We ran three independent
490 Markov chain Monte Carlo (MCMC) chains, each spanning at least 100 million iterations, and
491 sampled every 10,000 steps. The prior mutation rate derived from Pearson and colleagues was
492 used. Visual inspection of the trace from each MCMC chain confirmed signal convergence, with

493 effective sampling size values > 200 for key parameters. Visualisation of results were performed
494 using the R package ggtree⁷¹ v.3.10.0 to generate credibility time-calibrated phylogeny for each
495 sub-lineage (**Supplementary Figures 3**).

496

497 **Ancestral state reconstruction analysis**

498 Ancestral trait reconstruction was conducted to discern the dissemination patterns of *B.*
499 *pseudomallei* among provinces in northeast Thailand, focusing on dominant lineages. Due to
500 varying number of isolates among provinces, the analysis excluded Mahasarakham, which had
501 a limited dataset (n = 11), resulting in the analysis of eight provinces: Buriram, Khon Kaen,
502 Mukdahan, Nakhon Phanom, Roi Et, Sisaket, Surin, and Udon Thani. This approach yielded 28
503 potential province-to-province transmission combinations. To mitigate sampling biases, we sub-
504 sampled the phylogeny of each dominant lineage to have an equal number of isolates per
505 province (n = 15 isolates) and permuted 1,000 times. Using the stochastic character mapping
506 function (*make.simmap*) from the R package phytools⁷² v.1.9.16, we conducted 100 simulations
507 (nsim = 100) to reconstruct the provincial origins at each node in the sub-sampled phylogeny
508 (1,000 phylogenies per lineage). This allowed us to quantify transition events (Markov jumps)
509 between province pairs and determine the cumulative branch length associated with each
510 province (Markov rewards). A Mann-Whitney U test, with Bonferroni correction for multiple
511 comparisons was applied to compare the transition event counts among provinces
512 (**Supplementary Figure 6**).

513

514 **Pan-genome analysis**

515 All the study genomes were annotated using Prokka⁷³ v.1.14.5, and further used in the pan-
516 genome analysis. Each genome has a median of 5,845 coding sequences (CDS) predicted onto
517 each genome with a range of 5,642 to 6,142 CDS per genome. Panaroo⁷⁴ v.1.3.3 was
518 employed to estimate the pan-genome with a sensitive option and a cut-off sequence identity of
519 92% derived from previous study¹⁵. The number of estimated genes falls within a comparable
520 range to previous studies from a single population^{29,53}.

521

522 **Identification of dominant lineage-specific genes**

523 We determined lineage-specific genes by assessing their prevalence within dominant lineage or
524 any of their sub-lineages, requiring a high occurrence (95%) within these specific groups while
525 maintaining a low presence in non-dominant lineages. To achieve this, three thresholds were
526 employed: strict (95 % occurrence in dominants vs 5% occurrence in non-dominants),

527 intermediate (95% in dominants vs 10% in non-dominants), and relaxed (95% in dominants vs
528 15% in non-dominants). Based on visual examination of gene distribution patterns (**Figure 3a,**
529 **Supplementary Figure 7**), the relaxed threshold was used to maximise the number of genes
530 included in subsequent analysis.

531

532 **Identification of Gene Ontology (GO: terms)**

533 Amino acid sequences of lineage-specific genes were submitted to InterPro database⁴⁶
534 (<https://www.ebi.ac.uk/interpro/>) which characterised the function of lineage-specific genes
535 based on biological processes, molecular functions, and cellular compartments (**Figure 3;**
536 **Supplementary data 3**).

537

538 **Transcriptomic analysis of dominant lineage-specific genes**

539 **Lineage 1 transcriptome analysis**

540 The analysis focused on an expression profile of a strain K96243, which serves as a
541 representative of lineage 1. Data was sourced from microarray experiment generated by Ooi *et*
542 *al.*⁴⁷ and accessed through the Gene Expression Omnibus (GEO) under accession number
543 GSE43205.

544 To understand the difference between environmental and infection conditions, we compared the
545 expression profile of K96243 being exposed to water and K96243 recovered from infected mice.
546 To simulate environmental conditions, K96243 was cultured to log phase in LB medium,
547 subsequently washed with sterile deionising water, and suspended in water. To emulate
548 infection conditions, BALB mice were infected with 1000 CFU of *B. pseudomallei*, and bacteria
549 were harvested from the lungs three days post-infection. Two replicates were performed for
550 each condition. We retrieved microarray data from GEO using the R package GEOquery⁷⁵
551 v.2.58.0, and differential gene expression analysis was performed using the R package limma⁷⁶
552 v.3.58.1.

553 We used binary expression patterns reported in Ooi *et al.*⁴⁷ to compare the expression profiles
554 of lineage-specific genes against the remaining genes in K96243 across 62 conditions. This
555 enabled the comparison of the count of expressed genes within the lineage-specific category
556 against the remainder of the genes for each condition using Fisher's exact test, with multiple
557 testing adjustments via Benjamini-Hochberg corrections.

558 **Lineage 2 transcriptome analysis**

559 This analysis focused on the expression profile of a strain UKMD286, representative of lineage
560 2. RNAseq data was obtained from an experiment conducted by Ghazali *et al.*⁴⁸ and accessed
561 through the European Nucleotide Archive (ENA) (E-MTAB-11200).

562 To simulate environmental conditions, UKMD286 was cultured in BHIB medium overnight,
563 resuspended, and inoculated into soil extract medium. For infection condition, BALB mice were
564 infected with UKMD286, and the bacteria were harvested from spleens five days post infection.
565 Each experiment was conducted with two replicates. FastQC v.0.11.9 and FastXtool v.0.0.14
566 were used to pre-processed sequenced reads. Raw reads were aligned to UKDM286 genome
567 using Hisat2⁷⁷ v. 2.2.1 with differential gene expression performed using the R package
568 DESeq2⁷⁸ v.1.40.2.

569 **Lineage 3 transcriptome analysis**

570 We used a strain UKMH10 to represent the expression profile of lineage 3. Data was originated
571 from an RNAseq experiment conducted by Kong *et al.*⁴⁹ and was accessed through the
572 European Nucleotide Archive (ENA) (PRJEB53338)

573 To replicate environmental conditions, UKMH10 was cultured in LB medium overnight and sub-
574 cultured into soil extract medium. To simulate infection conditions, UKMH10 was cultured in LB
575 medium overnight and inoculated into human plasma, then incubated at 37 C to mimic the
576 human body temperature. Bacterial cells were harvested once the absorbance reading of the
577 bacterial cultures at 600 nm (OD600) reached 0.5. Each experiment was performed with two
578 replicates. FastQC v.0.11.9 and FastXtool v.0.0.14 were used to pre-processed sequenced
579 reads. Raw reads were aligned to UKMH10 genome using Hisat2⁷⁷ v.2.2.1 with differential gene
580 expression performed using the R package DESeq2⁷⁸ v.1.40.2.

581 **Acknowledgements**

582 We would like to thank physician and microbiology staff at the Udon Thani Hospital, Khon Kaen
583 Hospital, Srinakarin Hospital, Nakhon Phanom Hospital, Mukdahan Hospital, Roi Et Hospital,
584 Surin Hospital, Sisaket Hospital and Buriram Hospital for reporting the case and providing the
585 facilities for bacterial collection. We sincerely appreciate the transcriptomic data of *B.*
586 *pseudomallei* UKMD286 and UKMH10 provided by Professor Sheila Nathan and colleagues.
587 We are grateful to Dr. Soe Htet Aung for supporting geographical map plot. This work was
588 funded by Mahidol University (MU-KMUTT Biomedical engineering and Biomaterials

589 Consortium) and Royal Golden Jubilee Ph.D. Programme (RGJ-ASEAN) (<http://rgj.trf.or.th>)
590 through RS and NC. CCho was funded by Wellcome International Master Fellowship
591 (221418/Z/20/Z). CChe was funded by Wellcome International Intermediate Fellowship
592 (216457/Z/19/Z) and Sanger International Fellowship. NC and TEW were supported by the
593 National Institute of Allergy and Infectious Diseases of the National Institutes of Health
594 (NIAID/NIH) (<https://www.nih.gov>) under Award Number U01AI115520. The content is solely the
595 responsibility of the authors and does not necessarily represent the official views of the National
596 Institutes of Health. This research was funded in part, by the Wellcome Trust [220211]. For the
597 purpose of Open Access, the author has applied a CC BY public copyright license to any Author
598 Accepted Manuscript version arising from this submission. The funders had no role in study
599 design, data collection and analysis, decision to publish, or preparation of the manuscript.

600 **Data availability**

601 The genome sequence data presented in this study can be found in online repositories. The
602 ENA under study accession number PRJEB25606 and PRJEB35787. The accession numbers
603 for individual genomes, and their annotated assemblies are listed in **Supplementary data 1**.

604 **Code availability**

605 The analyses used public software including Kraken v.1.1.1, CheckM v.1.2.2, FastANI v.1.31,
606 FastQC v.0.11.9, FastXtool v.0.0.14, Unicycler v.0.8.4, ABACAS v.1.3.1, Velvet v.1.2.10,
607 Prokka v.1.13.4, Panaroo v.1.2.9, Snippy v.4.6.0, PopPUNK v.2.6.0, IQ-TREE v.2.0.3, Gubbins
608 v.3.1.3, Hisat2 v.2.2.1 and R packages: treespace v.1.1.4.3, BactDating v.1.1.1, ggtree v.3.10.0,
609 phytools v.1.9.16, GEOquery v.2.58.0, limma v.3.58.1 and DESeq2 v.1.40.2

610 **Ethics and Inclusion statements**

611 The research is led by local researchers and include local contributions throughout all research
612 process including the study design, study implementation, data ownership, intellectual
613 properties and authorship for publication.

614 **Author contributions**

615 RS, CChe and NC conceived, designed the study and write the original draft. NC and CChe
616 administrated and supervised the project. NC acquired funding to collect isolates, while CChe
617 acquired funding for sequencing and downstream analysis. NC, RS, TEW, and NS collected

618 and identified bacterial isolates. NC, RS and RP collected clinical data. CChе, NRT, NC, RS,
619 JT, EB and WC performed whole-genome sequencing. NPJD and NC contributed reagents.
620 NRT and JP contributed software tools. RS, CCho, and CChе performed bioinformatics
621 analyses. CChе, NC, NRT and JP interpret the analyses. CChе wrote the revised draft. All
622 authors read and approved the manuscript.

623 **Table 1** Recombination in dominant lineages.

Lineage	Percent of CDS impacted by recombination at least once (recombinant CDS/total CDSs in the reference genome)	Percent of dominant lineage-specific genes underwent recombination at least once (recombinant gene/lineage-specific genes identified in the reference genome)	Average r/m (number of SNPs introduced by recombination/ SNPs introduced by substitutions)		
			Internal nodes (95% CI)	Terminal nodes (95% CI)	Average (95% CI)
1	99.5% (5981/6010)	100% (47/47)	2.8 (2.3-3.3)	4.6 (4.0-5.3)	3.7 (3.3-4.1)
2	99.9% (5963/5971)	100% (65/65)	5.0 (3.9-6.0)	4.3 (3.6-4.9)	4.6 (4.0-5.2)
3	96.6% (5338/5523)	100% (20/20)	1.4 (1.1-1.8)	2.9 (2.2-3.7)	2.2 (1.8-2.6)

624

625 **Figure and table legend**

626 **Figure 1** Distribution of *B. pseudomallei* genomes used in this study (a) Geographical
627 representation of the countries and provinces sampled for the 1,391 *B. pseudomallei* genomes
628 used in this study. Pie-chart summarises the proportion of dominant lineage 1, 2, and 3
629 presented at each location with the chart size proportional to the number of the samples
630 collected (b) An unrooted phylogenetic tree colour-coded by dominant lineages (c) Histogram
631 depicting the distribution of clinical *B. pseudomallei* isolates from the northeast Thailand cohort
632 throughout 2015-2018 sampling period. The shaded blue area represents the period of rainy
633 seasons. (d) Boxplots summarising the pairwised core genome SNP distances among isolates
634 in this study, shown in a logarithmic scale. The distribution is depicted for the entire population
635 and each dominant lineage.

636 **Figure 2** Dissemination patterns in northeast Thailand. (a) Province-to-province transmission
637 patterns influenced by northeast Thailand geographical landscape. Nodes present provinces,
638 denoted by abbreviation and ordered by altitude: U-Udon Thani, K-Khon Kaen, B-Buriram, R -
639 Roi Et, N-Nakhon Phanom, Su-Surin, M-Mukdahan, and Si-Sisaket. Rivers are depicted in blue
640 with major rivers including the Great Mekong River, the Chi River, and the Mun River and their
641 flow direction annotated. (b) Average altitude of provinces in meters above sea level. Error bars
642 present 95% confidence interval. The northwest provinces exhibit higher altitudes, gradually
643 declining towards the southeast. For (a) and (b) solid arrows illustrating transmission
644 directionality explained by altitude differences. Dotted arrows represent transmission
645 directionality influenced by northeast monsoon winds. Grey arrows signify patterns with unclear
646 explanation.

647 **Figure 3** Dominant lineage-specific genes and their Gene Ontology (GO terms). (a) The
648 heatmap represents lineage-specific genes (right) detected in each isolate, aligned with the
649 phylogeny (left). Lineage-specific genes shared across multiple dominant lineages are
650 highlighted in yellow. Lineage-specific genes from lineage 1, 2, 3 are coloured in green, red, and
651 purple, respectively. Additionally, the colour stripes provide information on the lineage and sub-
652 lineage membership (b) Bar plots displays the frequency of GO annotations of lineage-specific
653 genes in each dominant lineage categorised by biological process, molecular function, and
654 cellular compartment. The pie-charts summarise the proportion of lineage-specific genes with
655 assigned GO terms (black).

656 **Figure 4** Transcriptome analysis of representative strains: K96243 (lineage 1), UKMD286
657 (lineage 2) and UKMH10 (lineage 3) (a to c) Volcano plots demonstrate differential gene
658 expression (DGE) between environmental and infection conditions. Vertical dotted lines
659 represent the statistical cut-off at log two-fold change, while horizontal dotted lines display the
660 statistical cut-off at the adjusted p-value of 0.05 on a negative log scale. Each dot represents a
661 gene, with lineage-1, lineage-2, and lineage-3-specific coloured in green, red, and purple,
662 respectively. (d) Binary expression profile of lineage-1-specific genes across different
663 conditions. A star denotes significant differences in the gene expression profile of lineage-1-
664 specific genes compared the remaining genes of strain K96243.

665 **Supplementary Figure 1** Comparison of approaches used in in outlining the population
666 structure (a) A scatter plot displays the first two principal components (PCs) derived from
667 pairwise distances between trees. Each dot represents a bootstrap tree and is colour-coded by
668 the method used to generate the tree: core genome SNP (blue), cgMLST (red) and MLST
669 (green). (b) A scree plot summarises eigenvalues computed for each PCs. (c to d) The median
670 phylogenetic trees constructed from core genome SNP, cgMLST, and MLST and their
671 consistency with PopPUNK clustering method.

672 **Supplementary Figure 2** Distribution of environmental and clinical isolates by each lineage.
673 Barplots highlight the co-detection of environmental (green) and clinical isolates (red) across
674 dominant lineages 1, 2, and 3.

675 **Supplementary Figure 3** Lineage 1 specific analysis (a) A recombination removed lineage 1
676 phylogeny with colour stripes displaying its sub-lineage structure, year of collection, and
677 sampling province (left to right). (b) A map of northeast Thailand showing the distribution of
678 each isolate and the region's river system. (c) Time-calibrated phylogeny of sub-lineage 1.3 with
679 blue error bars indicating 95% highest posterior density interval, with the estimated mutational
680 rate consistent with previous study.

681 **Supplementary Figure 4** Lineage 2 specific analysis (a) A recombination removed lineage 2
682 phylogeny with colour stripes highlighting sub-lineage structure, year of collection, and sampling
683 province (left to right). (b) A map of northeast Thailand with the region's river system. Dots
684 present the distribution of individual samples.

685 **Supplementary Figure 5** Lineage 3 specific analysis (a) A recombination removed lineage 3
686 phylogeny with colour stripes highlighting sub-lineage structure, year of collection, and sampling

687 province (left to right). (b) A map of northeast Thailand showing the distribution of each isolate
688 and the region's river system.

689 **Supplementary Figure 6** Transmission patterns and evolutionary time spent at each province.
690 (a to c) Proportion of transition events (Markov jumps) among provincial pairs for lineage 1, 2,
691 and 3 respectively. The pairs were denoted as province1 – province 2, with transitions from
692 province 1 to province 2 shown in red, and transitions from province 2 to province 1 in green.
693 The Man-Whitney U test was conducted for each pair to assess differences in transition
694 frequency by direction, with Bonferroni correction applied for multiple tests. (d to f) Total branch
695 length from provincial trait reconstruction (Markov rewards) for lineage 1, 2, and 3, respectively.

696 **Supplementary Figure 7** Selection criteria for lineage-specific genes. Scatter plots show the
697 frequency distribution of lineage-specific (red) against other genes (black), based on their
698 distribution within the dominant lineages and their sub-lineages (horizontal axis) compared to
699 their distribution in non-dominant lineages.

700 **Supplementary Figure 8** Recombination patterns detected in lineage 1, 2, and 3. From left to
701 right: the recombination-removed phylogeny of each lineage, a stripe representing the sampling
702 year and sub-lineage classification, and heatmaps displaying recombination patterns identified
703 in chromosome 1 and 2. The top orange lines mark the genome coordinates. For each lineage,
704 their respective lineage-specific genes are highlighted in blue at the top of the panel. Each
705 heatmap represents recombination blocks aligned with the phylogeny. Recombination events
706 occurring at the internal nodes are coloured in red, while those occurring at the external
707 branches are coloured in blue. The recombination hotspot is plotted at the bottom of each
708 heatmap.

709 **Supplementary data 1** Epidemiological data, isolate and accession codes for both short reads
710 and annotated assembly used in this study deposited in the European Nucleotide Archive (ENA)
711 (n = 1,391) (provided as a separate excel file).

712

713 **Supplementary data 2** Quality control information of newly sequenced *B. pseudomallei* from
714 northeast Thailand collection (n = 1,265) (provided as a separate excel file).

715

716 **Supplementary data 3** List of dominant lineage-specific genes and their Gene Ontology (GO)
717 terms (provided as a separate excel file).

718

719 **Supplementary data 4** List of dominant lineage-specific genes and their expression profile
720 under infectious and environmental conditions generated by Ooi *et al.* 2013, Ghazali *et al.* 2023
721 and Kong *et al.* 2023 (provided as a separate excel file).

722 **References**

- 723 1. Limmathurotsakul, D. *et al.* Predicted global distribution of *Burkholderia pseudomallei* and
724 burden of melioidosis. *Nat Microbiol* **1**, 15008 (2016).
- 725 2. Birnie, E. *et al.* Global burden of melioidosis in 2015: a systematic review and data synthesis.
726 *The Lancet Infectious Diseases* **19**, 892–902 (2019).
- 727 3. Mohapatra, P. R. & Mishra, B. Burden of melioidosis in India and South Asia: Challenges and
728 ways forward. *The Lancet Regional Health - Southeast Asia* **2**, 100004 (2022).
- 729 4. Zheng, X., Xia, Q., Xia, L. & Li, W. Endemic Melioidosis in Southern China: Past and Present.
730 *TropicalMed* **4**, 39 (2019).
- 731 5. Tran, Q. T. L. *et al.* Child Melioidosis Deaths Caused by *Burkholderia pseudomallei* –
732 Contaminated Borehole Water, Vietnam, 2019. *Emerg. Infect. Dis.* **28**, 1689–1693 (2022).
- 733 6. Win, M. M. *et al.* Enhanced melioidosis surveillance in patients attending four tertiary
734 hospitals in Yangon, Myanmar. *Epidemiol. Infect.* **149**, e154 (2021).
- 735 7. Thabit, A. M. *et al.* Etiologies of tropical acute febrile illness in West Pahang, Malaysia: A
736 prospective observational study. *Asian Pac J Trop Med* **13**, 115 (2020).
- 737 8. Chantratita, N. *et al.* Characteristics and one year outcomes of melioidosis patients in
738 Northeastern Thailand: A prospective, multicenter cohort study. *The Lancet Regional Health -*
739 *Southeast Asia* **9**, 100118 (2023).
- 740 9. Hodgetts, K. *et al.* Melioidosis in the remote Katherine region of northern Australia. *PLoS*
741 *Negl Trop Dis* **16**, e0010486 (2022).
- 742 10. Gassiep1, I., Ganeshalingam, V., Chatfield, M. D., Harris, P. N. A. & Norton, R. E. The
743 epidemiology of melioidosis in Townsville, Australia. *Transactions of The Royal Society of*
744 *Tropical Medicine and Hygiene* **116**, 328–335 (2022).
- 745 11. Hinjoy, S. *et al.* Melioidosis in Thailand: Present and Future. *TropicalMed* **3**, 38 (2018).
- 746 12. Bulterys, P. L. *et al.* Climatic drivers of melioidosis in Laos and Cambodia: a 16-year
747 case series analysis. *Lancet Planet Health* **2**, e334–e343 (2018).
- 748 13. Chai, L. Y. A. & Fisher, D. Earth, wind, rain, and melioidosis. *The Lancet Planetary*
749 *Health* **2**, e329–e330 (2018).
- 750 14. Liu, X. *et al.* Association of Melioidosis Incidence with Rainfall and Humidity, Singapore,
751 2003–2012. *Emerg. Infect. Dis.* **21**, 159–162 (2015).

- 752 15. Chewapreecha, C. *et al.* Global and regional dissemination and evolution of
753 *Burkholderia pseudomallei*. *Nat Microbiol* **2**, 16263 (2017).
- 754 16. Zheng, H. *et al.* Genetic diversity and transmission patterns of *Burkholderia*
755 *pseudomallei* on Hainan island, China, revealed by a population genomics analysis. *Microbial*
756 *Genomics* **7**, (2021).
- 757 17. Sarovich, D. S. *et al.* Phylogenomic Analysis Reveals an Asian Origin for African
758 *Burkholderia pseudomallei* and Further Supports Melioidosis Endemicity in Africa. *mSphere*
759 **1**, e00089-15 (2016).
- 760 18. Price, E. P., Currie, B. J. & Sarovich, D. S. Genomic Insights Into the Melioidosis
761 Pathogen, *Burkholderia pseudomallei*. *Curr Trop Med Rep* **4**, 95–102 (2017).
- 762 19. Mayo, M. *et al.* A cluster of melioidosis cases from an endemic region is clonal and is
763 linked to the water supply using molecular typing of *Burkholderia pseudomallei* isolates. *The*
764 *American Journal of Tropical Medicine and Hygiene* **65**, 177–179 (2001).
- 765 20. Draper, A. D. K. *et al.* Association of the Melioidosis Agent *Burkholderia pseudomallei*
766 with Water Parameters in Rural Water Supplies in Northern Australia. *Appl Environ Microbiol*
767 **76**, 5305–5307 (2010).
- 768 21. Chen, Y.-L. *et al.* The Concentrations of Ambient *Burkholderia Pseudomallei* during
769 Typhoon Season in Endemic Area of Melioidosis in Taiwan. *PLoS Negl Trop Dis* **8**, e2877
770 (2014).
- 771 22. Chen, P.-S. *et al.* Airborne Transmission of Melioidosis to Humans from Environmental
772 Aerosols Contaminated with *B. pseudomallei*. *PLoS Negl Trop Dis* **9**, e0003834 (2015).
- 773 23. Limmathurotsakul, D. *et al.* Melioidosis Caused by *Burkholderia pseudomallei* in Drinking
774 Water, Thailand, 2012. *Emerg. Infect. Dis.* **20**, 265–268 (2014).
- 775 24. Yip, T.-W. *et al.* Endemic Melioidosis in Residents of Desert Region after Atypically
776 Intense Rainfall in Central Australia, 2011. *Emerg. Infect. Dis.* **21**, 1038–1040 (2015).
- 777 25. Baker, A. L., Ezzahir, J., Gardiner, C., Shipton, W. & Warner, J. M. Environmental
778 Attributes Influencing the Distribution of *Burkholderia pseudomallei* in Northern Australia.
779 *PLoS ONE* **10**, e0138953 (2015).
- 780 26. Pumpuang, A. *et al.* Survival of *Burkholderia pseudomallei* in distilled water for 16 years.
781 *Transactions of the Royal Society of Tropical Medicine and Hygiene* **105**, 598–600 (2011).
- 782 27. Chewapreecha, C. *et al.* Co-evolutionary Signals Identify *Burkholderia pseudomallei*
783 Survival Strategies in a Hostile Environment. *Mol Biol Evol* **39**, msab306 (2022).
- 784 28. Hantrakun, V. *et al.* Soil Nutrient Depletion Is Associated with the Presence of
785 *Burkholderia pseudomallei*. *Appl. Environ. Microbiol.* **82**, 7086–7092 (2016).

- 786 29. Chewapreecha, C. *et al.* Genetic variation associated with infection and the environment
787 in the accidental pathogen *Burkholderia pseudomallei*. *Commun Biol* **2**, 428 (2019).
- 788 30. Holden, M. T. G. *et al.* Genomic plasticity of the causative agent of melioidosis,
789 *Burkholderia pseudomallei*. *Proceedings of the National Academy of Sciences* **101**, 14240–
790 14245 (2004).
- 791 31. Nandi, T. *et al.* *Burkholderia pseudomallei* sequencing identifies genomic clades with
792 distinct recombination, accessory, and epigenetic profiles. *Genome Res.* **25**, 129–141 (2015).
- 793 32. Liechti, N. *et al.* Whole-Genome Assemblies of 16 *Burkholderia pseudomallei* Isolates
794 from Rivers in Laos. *Microbiol Resour Announc* **10**, e01226-20 (2021).
- 795 33. Rachlin, A. *et al.* Whole-genome sequencing of *Burkholderia pseudomallei* from an
796 urban melioidosis hot spot reveals a fine-scale population structure and localised spatial
797 clustering in the environment. *Sci Rep* **10**, 5443 (2020).
- 798 34. Webb, J. R. *et al.* Myanmar *Burkholderia pseudomallei* strains are genetically diverse
799 and originate from Asia with phylogenetic evidence of reintroductions from neighbouring
800 countries. *Sci Rep* **10**, 16260 (2020).
- 801 35. Lees, J. A. *et al.* Fast and flexible bacterial genomic epidemiology with PopPUNK.
802 *Genome Res.* **29**, 304–316 (2019).
- 803 36. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic
804 Inference in the Genomic Era. *Molecular Biology and Evolution* **37**, 1530–1534 (2020).
- 805 37. Lichtenegger, S. *et al.* Development and Validation of a *Burkholderia pseudomallei* Core
806 Genome Multilocus Sequence Typing Scheme To Facilitate Molecular Surveillance. *J Clin*
807 *Microbiol* **59**, e00093-21 (2021).
- 808 38. Godoy, D. *et al.* Multilocus Sequence Typing and Evolutionary Relationships among the
809 Causative Agents of Melioidosis and Glanders, *Burkholderia pseudomallei* and *Burkholderia*
810 *mallei*. *J Clin Microbiol* **41**, 2068–2079 (2003).
- 811 39. Lofjle, E. & Kubiniok, J. Landform Development and Bioturbation of the Khorat Plateau,
812 Northeast Thailand. *Nat Hist Bull Siam Soc* 199–216 (1996).
- 813 40. NASA. Thailand Elevation Map.
814 <https://www.floodmap.net/elevation/CountryElevationMap/?ct=TH> (2020).
- 815 41. Thai Meteorological Department. The Climate of Thailand. (2016).
- 816 42. Arpornrat, T., Ratjiranukool, S., Ratjiranukool, P. & Sasaki, H. Evaluation of Southwest
817 Monsoon Change over Thailand by High-resolution Regional Climate Model under High RCP
818 Emission Scenario. *J. Phys.: Conf. Ser.* **1144**, 012112 (2018).

- 819 43. Hien, P. D., Bac, V. T. & Thinh, N. T. H. PMF receptor modelling of fine and coarse
820 PM10 in air masses governing monsoon conditions in Hanoi, northern Vietnam. *Atmospheric*
821 *Environment* **38**, 189–201 (2004).
- 822 44. Hien, T. T. *et al.* Current Status of Fine Particulate Matter (PM2.5) in Vietnam's Most
823 Populous City, Ho Chi Minh City. *Aerosol Air Qual. Res.* **19**, 2239–2251 (2019).
- 824 45. Floch, P. & Molle, F. Marshalling water resources: A chronology of irrigation
825 development in the Chi-Mun River basin Northeast Thailand. (2007).
- 826 46. The Gene Ontology Consortium. The Gene Ontology Resource: 20 years and still GOing
827 strong. *Nucleic Acids Research* **47**, D330–D338 (2019).
- 828 47. Ooi, W. F. *et al.* The Condition-Dependent Transcriptional Landscape of Burkholderia
829 pseudomallei. *PLoS Genet* **9**, e1003795 (2013).
- 830 48. Ghazali, A.-K. *et al.* Transitioning from Soil to Host: Comparative Transcriptome Analysis
831 Reveals the Burkholderia pseudomallei Response to Different Niches. *Microbiol Spectr* **11**,
832 e03835-22 (2023).
- 833 49. Kong, C. *et al.* Transcriptional landscape of Burkholderia pseudomallei cultured under
834 environmental and clinical conditions. *Microbial Genomics* **9**, (2023).
- 835 50. Manivanh, L. *et al.* Burkholderia pseudomallei in a lowland rice paddy: seasonal
836 changes and influence of soil depth and physico-chemical properties. *Sci Rep* **7**, 3031
837 (2017).
- 838 51. Tuanyok, A. *et al.* Genomic islands from five strains of Burkholderia pseudomallei. *BMC*
839 *Genomics* **9**, 566 (2008).
- 840 52. Bertelli, C. *et al.* Enabling genomic island prediction and comparison in multiple
841 genomes to investigate bacterial evolution and outbreaks. *Microbial Genomics* **8**, (2022).
- 842 53. Spring-Pearson, S. M. *et al.* Pangenome Analysis of Burkholderia pseudomallei:
843 Genome Evolution Preserves Gene Order despite High Recombination Rates. *PLoS ONE*
844 **10**, e0140274 (2015).
- 845 54. *White gold: the commercialisation of rice farming in the Lower Mekong Basin.* (Palgrave
846 Macmillan, 2020).
- 847 55. Keyes, C. F. *Isan: Regionalism in Northeastern Thailand.* (Southeast Asia Program,
848 Department of Asian Studies, Cornell University, 1967).
- 849 56. Charusiri, P., Daorerk, V., Archibald, D., Hisada, K. & Ampaiwan, T. Geotectonic
850 Evolution of Thailand: A new synthesis. *Journal of the Geological Society of Thailand* 1–20
851 (2002).

- 852 57. Borrelli, P. *et al.* An assessment of the global impact of 21st century land use change on
853 soil erosion. *Nat Commun* **8**, 2013 (2017).
- 854 58. Singh, B. K. *et al.* Climate change impacts on plant pathogens, food security and paths
855 forward. *Nat Rev Microbiol* **21**, 640–656 (2023).
- 856 59. Wood, D. E. & Salzberg, S. L. Kraken: ultrafast metagenomic sequence classification
857 using exact alignments. *Genome Biol* **15**, R46 (2014).
- 858 60. Parks, D. H., Imelfort, M., Skennerton, C. T., Hugenholtz, P. & Tyson, G. W. CheckM:
859 assessing the quality of microbial genomes recovered from isolates, single cells, and
860 metagenomes. *Genome Res.* **25**, 1043–1055 (2015).
- 861 61. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High
862 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat*
863 *Commun* **9**, 5114 (2018).
- 864 62. Zerbino, D. R. & Birney, E. Velvet: Algorithms for de novo short read assembly using de
865 Bruijn graphs. *Genome Res.* **18**, 821–829 (2008).
- 866 63. Page, A. J. *et al.* SNP-sites: rapid efficient extraction of SNPs from multi-FASTA
867 alignments. *Microbial Genomics* **2**, (2016).
- 868 64. Jombart, T., Kendall, M., Almagro-Garcia, J. & Colijn, C. TREESPACE²: Statistical
869 exploration of landscapes of phylogenetic trees. *Molecular Ecology Resources* **17**, 1385–
870 1392 (2017).
- 871 65. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic
872 tree display and annotation. *Nucleic Acids Research* **49**, W293–W296 (2021).
- 873 66. Wick, R. R., Judd, L. M., Gorrie, C. L. & Holt, K. E. Unicycler: Resolving bacterial
874 genome assemblies from short and long sequencing reads. *PLoS Comput Biol* **13**, e1005595
875 (2017).
- 876 67. Assefa, S., Keane, T. M., Otto, T. D., Newbold, C. & Berriman, M. ABACAS: algorithm-
877 based automatic contiguation of assembled sequences. *Bioinformatics* **25**, 1968–1969
878 (2009).
- 879 68. Croucher, N. J. *et al.* Rapid phylogenetic analysis of large samples of recombinant
880 bacterial whole genome sequences using Gubbins. *Nucleic Acids Research* **43**, e15–e15
881 (2015).
- 882 69. Tonkin-Hill, G., Lees, J. A., Bentley, S. D., Frost, S. D. W. & Corander, J. Fast
883 hierarchical Bayesian analysis of population structure. *Nucleic Acids Research* **47**, 5539–
884 5549 (2019).

- 885 70. Didelot, X., Croucher, N. J., Bentley, S. D., Harris, S. R. & Wilson, D. J. Bayesian
886 inference of ancestral dates on bacterial phylogenetic trees. *Nucleic Acids Research* **46**,
887 e134–e134 (2018).
- 888 71. Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. GGTREE: an R package for
889 visualization and annotation of phylogenetic trees with their covariates and other associated
890 data. *Methods Ecol Evol* **8**, 28–36 (2017).
- 891 72. Revell, L. J. phytools: an R package for phylogenetic comparative biology (and other
892 things). *Methods Ecol Evol* **3**, 217–223 (2012).
- 893 73. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–
894 2069 (2014).
- 895 74. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo
896 pipeline. *Genome Biol* **21**, 180 (2020).
- 897 75. Davis, S. & Meltzer, P. S. GEOquery: a bridge between the Gene Expression Omnibus
898 (GEO) and BioConductor. *Bioinformatics* **23**, 1846–1847 (2007).
- 899 76. Ritchie, M. E. *et al.* limma powers differential expression analyses for RNA-sequencing
900 and microarray studies. *Nucleic Acids Research* **43**, e47–e47 (2015).
- 901 77. Kim, D., Paggi, J. M., Park, C., Bennett, C. & Salzberg, S. L. Graph-based genome
902 alignment and genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* **37**, 907–915
903 (2019).
- 904 78. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion
905 for RNA-seq data with DESeq2. *Genome Biol* **15**, 550 (2014).
- 906

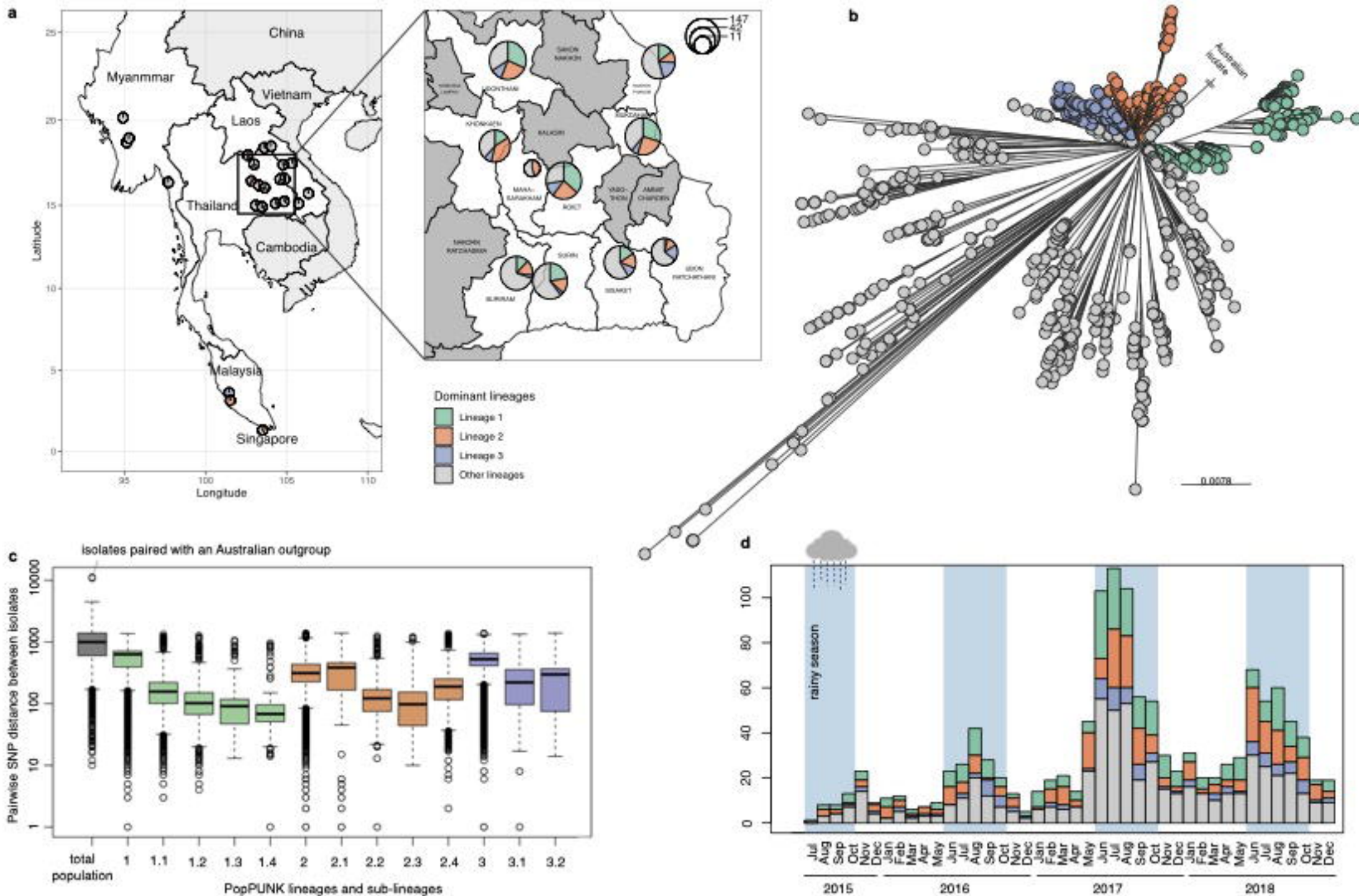
Figure 1

Figure 2

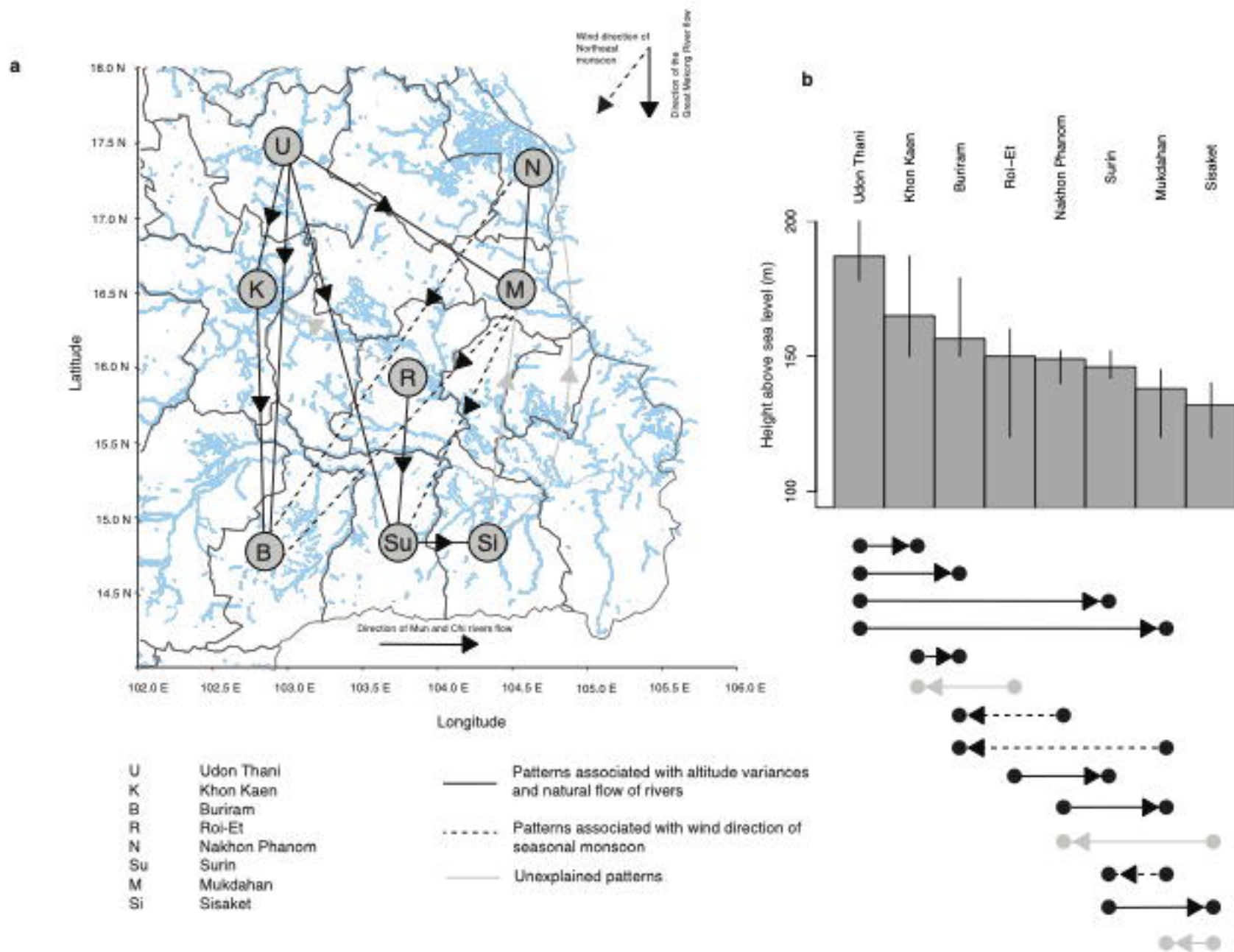


Figure 3

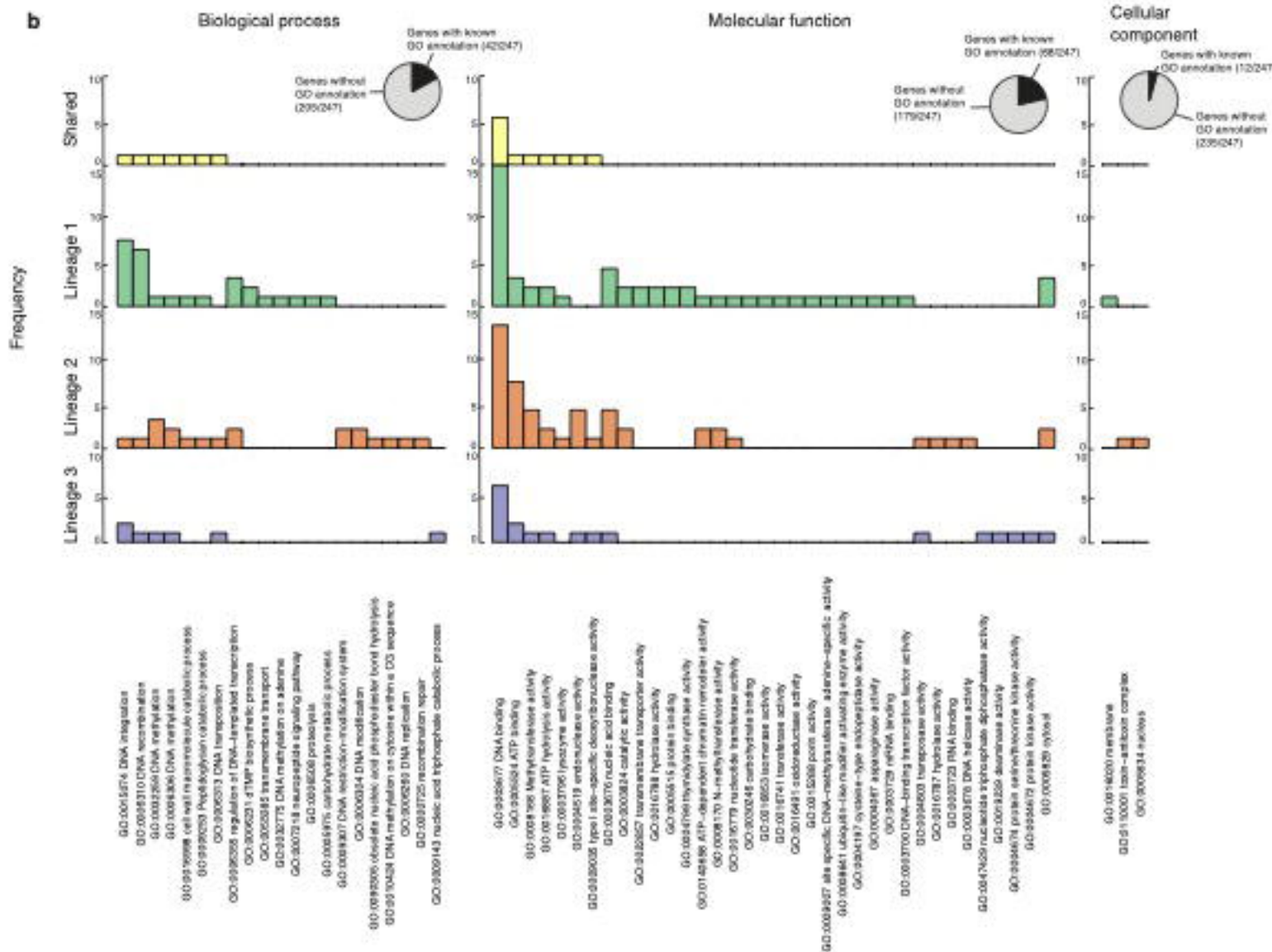
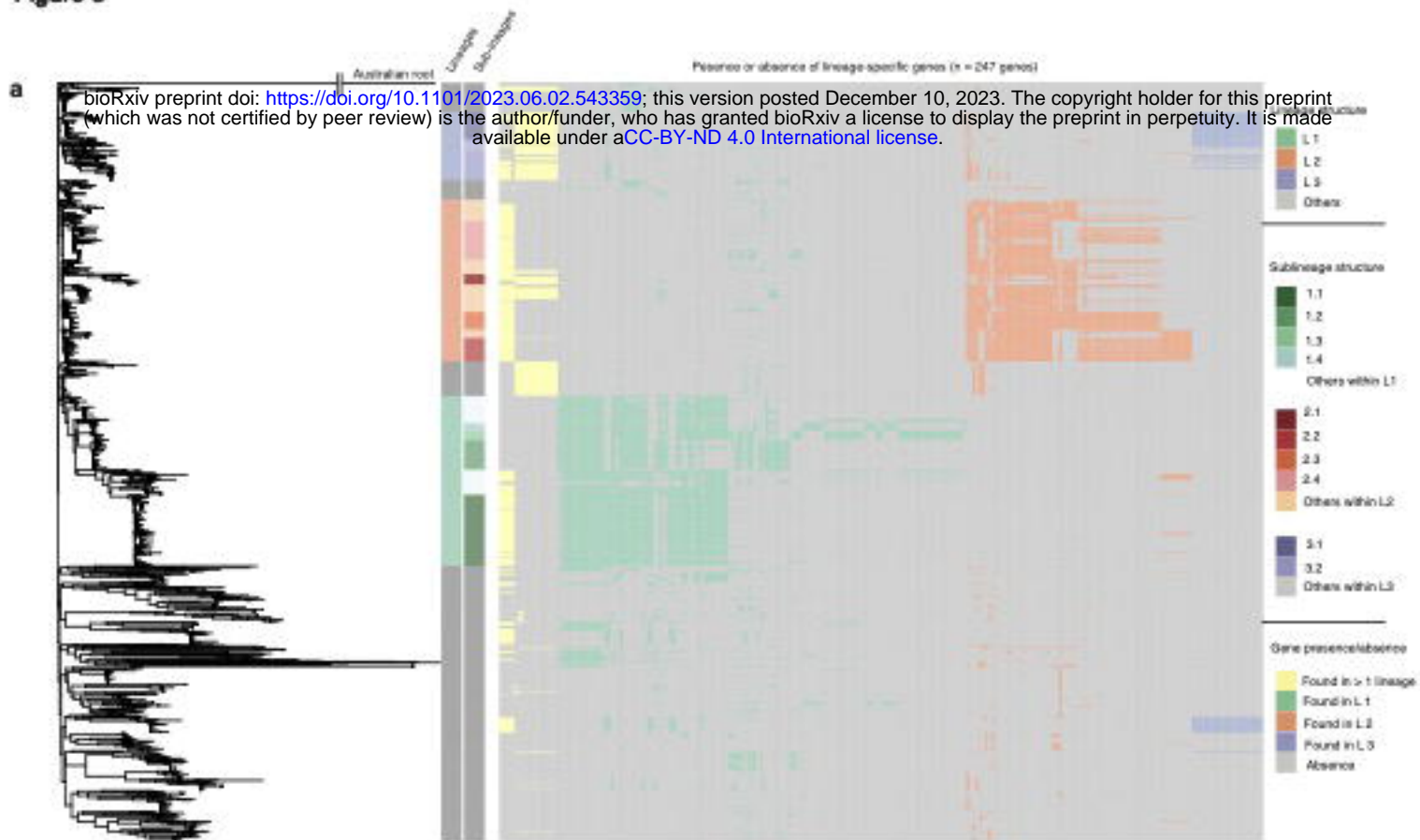


Figure 4

