

1 **Insights into the genetic diversity of *Mycobacterium***

2 ***tuberculosis* in Tanzania**

3 **Liliana K. Rutaihwa^{1, 2, 3*}, Mohamed Sasamalo^{1, 2,3}, Aladino Jaleco^{1, 2}, Jerry Hella^{1, 2,3}, Ally**
4 **Kingazi³, Lujeko Kamwela^{1, 2,3}, Amri Kingalu^{4, 5}, Bryceson Malewo^{4, 5}, Raymond Shirima^{4, 5},**
5 **Anna Doetsch^{1, 2}, Julia Feldmann^{1, 2}, Miriam Reinhard^{1, 2}, Sonia Borrell^{1, 2}, Klaus Reither^{1, 2},**
6 **Basra Doulla^{4, 5}, Lukas Fenner^{1, 2, 6#} and Sebastien Gagneux^{1, 2, *, #}**

7

8

9 ¹ Swiss Tropical and Public Health Institute, Basel, Switzerland

10 ² University of Basel, Basel, Switzerland

11 ³ Ifakara Health Institute, Bagamoyo, Tanzania

12 ⁴ Central Tuberculosis Reference Laboratory, Dar es Salaam, Tanzania

13 ⁵ National Tuberculosis and Leprosy Programme, Dar es Salaam, Tanzania

14 ⁶ Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland

15

16

17 * Corresponding authors

18 Email: liliana.rutaihwa@gmail.com (LKR) and sebastien.gagneux@swisstph.ch (SG)

19

20

21 # Equal contribution

22

23 Abstract

24 **Background:** Human tuberculosis (TB) is caused by seven phylogenetic lineages of the
25 *Mycobacterium tuberculosis* complex (MTBC), Lineage 1–7. Recent advances in rapid
26 genotyping of MTBC based on single nucleotide polymorphisms (SNP), allow for rapid and
27 phylogenetically robust strain classification, paving the way for defining genotype-phenotype
28 relationships in clinical settings. Such studies have revealed that, in addition to host and
29 environmental factors, different strains of the MTBC influence the outcome of TB infection
30 and disease. In Tanzania, such molecular epidemiological studies of TB however are scarce in
31 spite of a high TB burden.

32 **Methods and Findings:** Here we used a SNP-typing method to genotype a nationwide
33 collection of 2,039 MTBC clinical isolates obtained from new and retreatment TB cases
34 diagnosed in 2012 and 2013. Four lineages, namely Lineage 1–4 were identified. The
35 distribution and frequency of these lineages varied across the regions but overall, Lineage 4
36 was the most frequent (n = 866, 42.5%), followed by Lineage 3 (n = 681, 33.4%) and 1 (n =
37 336, 16.5%), with Lineage 2 being the least frequent (n = 92, 4.5%). A total of 64 (3.1%)
38 isolates could not be assigned to any lineage. We found Lineage 2 to be associated with
39 female sex (adjusted odds ratio [aOR] 2.25; 95% confidence interval [95% CI] 1.38 – 3.70, p <
40 0.001) and retreatment (aOR 1.78; 95% CI 1.00 – 3.02, p = 0.040). We found no associations
41 between MTBC lineage and patient age or HIV status. Our sublineage typing based on spacer
42 oligotyping revealed the presence of mainly EAI, CAS and LAM families. Finally, we detected
43 low levels of multidrug resistant isolates among a subset of retreatment cases

44 **Conclusions:** This study provides novel insights into the influence of pathogen-related
45 factors on the TB epidemic in Tanzania.

46

47 Introduction

48 Tuberculosis (TB) is the leading cause of mortality due to an infectious disease [1]. In 2017,
49 an estimated 10.0 million people developed TB globally, with 1.3 million dying of the
50 disease. More than 80% of the TB burden lies in the thirty high burden countries [1].
51 Tanzania is among these countries, with a national average TB notification rate of 129 cases
52 per 100,000; however, some regions show higher notification rates [2]. Like in most sub-
53 Saharan African countries, the HIV epidemic contributes to the high TB incidence in
54 Tanzania, where a-third of the TB patients are co-infected with HIV [2]. Contrarily, drug
55 resistant-TB is still low in this setting [3]. Other risk factors such poverty also influence the
56 epidemiology of TB in Tanzania [4].

57 Transmission of TB occurs via infectious aerosols, and upon exposure individuals can either
58 develop active disease or remain latent infected [5]. It is estimated that a-quarter of the
59 world's population is latently TB infected [6], with a 5 – 10% life time risk to develop active
60 TB disease; this risk is 50% in case of HIV co-infected individuals [7].

61 The complex dynamics of TB infection and disease are determined by the environment, the
62 host and the pathogen [8]. Seven main phylogenetic lineages of the *Mycobacterium*
63 *tuberculosis* complex (MTBC) lineages (Lineage 1–7) are known to cause TB in humans [9].
64 These lineages are phylogeographically distributed, reflecting human migration histories and
65 possibly adaptation to different human populations [10–12]. Some genomic differences
66 among the MTBC strains translate into relevant biological and epidemiological phenotypes
67 [13]. In general, strains of the globally distributed lineages, Lineage 2 and 4 or “generalists”,
68 appear to be more virulent in average than those of the geographically restricted lineages,

69 Lineage 5 and 6 or “specialists” [9,13]. Epidemiologically speaking, these phenotypes are
70 demonstrated by indicators such as transmission potential, disease severity and rate of
71 progression from infection to disease [14–17].

72 Studying genotype-phenotype relationships requires understanding the genetic diversity of
73 MTBC clinical strains in a given clinical setting. In Tanzania few studies have described the
74 genetic diversity of MTBC [18–21]. These previous work revealed the presence of mainly
75 three lineages; Lineage 1, 3 and 4, which include the EAI, CAS and LAM spoligo families,
76 respectively. Lineage 2, which includes the Beijing family, has only been reported at the
77 lowest frequencies. These previous studies are limited as they only focused on few
78 geographical locations and used spacer oligonucleotide typing (spoligotyping) technique
79 which has limitations for phylogenetically robust strain classification [22,23]. Only one study
80 profiled MTBC on a countrywide scale albeit with low sampling coverage [20].

81 In this study we used for the first time a robust single nucleotide polymorphism (SNP) typing
82 method to classify the largest so far nationwide collection of MTBC clinical isolates from
83 Tanzania. We then looked for potential associations between the MTBC lineages and the
84 clinical and epidemiological characteristics of the patients.

85

86 **Material and Methods**

87 **Study setting**

88 Our study was based on a nationwide convenience sample of sputum smear positive new
89 and retreatment TB cases diagnosed between 2012 and 2013 in Tanzania. The collection was
90 obtained via a platform established for routine TB drug resistance surveillance by the
91 National Tuberculosis Leprosy Program (NTLP) of Tanzania, covering health facilities in all
92 geographical regions of the country. Briefly, smear positive sputa specimens from
93 approximately 25% of new TB cases (obtained by allocating four months of sample collection
94 to each region annually) and from all retreatment cases were sent to zonal reference
95 laboratories (i.e. Central Tuberculosis Reference Laboratory [CTRL] in Dar es Salaam,
96 Bugando Medical Center [BMC] in Mwanza and Kilimanjaro Christian Medical Center [KCMC]
97 in Kilimanjaro), which serve the respective nearby regions for culture. Isolates from the two
98 zonal laboratories, BMC and KCMC were then sent to the CTRL for drug susceptibility testing
99 (DST). For this study we used archived isolates obtained from the CTRL.

100 **Study population and data collection**

101 We included a total of 2,039 unique (single patient) culture-confirmed TB cases diagnosed
102 between 2012 and 2013, each of whom we could retrieve the respective culture isolate. This
103 study population represents 41% of all culture positive sputum samples processed and 1.6%
104 of all TB notified cases in the country during the study period (S1 Fig). We also obtained
105 corresponding socio-demographic and clinical information collected during patients'
106 consultation at the respective health facilities. The demographic data collected included age,
107 sex and geographical location of the patients, whereas clinical data included HIV status and
108 disease category (i.e., new case and retreatment case).

109 **Processing of culture isolates**

110 The smear positive sputa samples were cultured on Löwenstein Jensen (LJ) growth medium
111 according to laboratory protocols. For this study, we included *M. tuberculosis* clinical isolates
112 retrieved from archived LJ media. We then prepared heat inactivated samples for the
113 retrieved clinical isolates by suspending *M. tuberculosis* colonies into 1ml sterile water and
114 heat inactivate at 95°C for one hour.

115 **Molecular genotyping**

116 We then classified the *M. tuberculosis* clinical isolates into main phylogenetic lineages by
117 TaqMan real-time PCR according to standard protocols (Applied Biosystems, Carlsbad, USA)
118 and as previously described [24]. We also performed 43-spacer spoligotyping on a
119 membrane for a subset of representative *M. tuberculosis* clinical strains following standard
120 protocols [25]. The clinical strains were assigned to spoligotype families using the online
121 database SITVITWEB [26].

122 **Drug Resistance Genotyping**

123 We selected a subset of clinical isolates from retreatment cases to perform molecular drug
124 resistance testing. We used a previously described multiplex polymerase chain reaction
125 (PCR) to target the hotspot region of *rpoB* gene that confers resistance to rifampicin [27].
126 The PCR assay targets both the tuberculous and non-tuberculous *Mycobacteria* (MTBC and
127 NTMs, respectively) *rpoB* gene, so we could also rule out the presence of non-tuberculous
128 isolates in our study sample using the assay. The amplified *rpoB* gene product was confirmed
129 by electrophoresis on a 2% agarose gel and sent for Sanger sequencing. We analyzed the
130 resulting sequences by Staden software package [28] and using *M. tuberculosis* H37Rv *rpoB*
131 gene as reference sequence.

132 **Statistical Analyses**

133 For statistical analyses we applied descriptive statistics to delineate patients' characteristics.

134 We used Chi-square or Fisher's exact tests for assessment of differences between groups in

135 categorical variables, whenever applicable. We used univariate and multivariate logistic

136 regression models to assess for the association between *M. tuberculosis* lineages and

137 patients' clinical and demographic characteristics. The associations were assessed for

138 Lineage 2 compared to all other lineages (Lineages 1, 3 and 4), adjusting for age, sex, disease

139 category and HIV status. All statistical analyses were performed in R 3.5.0 [29].

140

141 **Results**

142 **Patients' demographic and clinical characteristics**

143 The patients' demographic and clinical information in our study included; age, sex,
144 geographical location, HIV and disease category (new or retreatment case). Table 1 describes
145 patients' characteristics of the study population. The proportions of the observed and
146 missing data for the study population are summarized in S2 Fig.

147

148 **Table 1.** Clinical and demographic characteristics of the TB cases

Characteristics	Valid Proportion %	Total (%) n = 2039
Age, median (IQR)		
28 (20-38)		
Age groups (years)		
Child age (< 15)	9.87	193 (9.47)
Young age (15 - 24)	29.67	580 (28.45)
Early adult (25 - 44)	47.98	938 (46.00)
Late adult (45 - 64)	10.03	196 (9.61)
Old age (> 65)	2.46	48 (2.35)
Not available		84 (4.12)
	<i>total n = 1955</i>	
Sex		
Female	32.40	645 (31.63)
Male	67.60	1346 (66.01)
Not available		48 (2.35)
	<i>total n = 1991</i>	
HIV status		
Negative	67.71	1086 (53.26)
Positive	32.23	517 (25.36)
Indeterminate	0.06	1 (0.05)
Not available		435 (21.33)
	<i>total n = 1604</i>	
Patient category		
New case	83.95	1679 (82.34)
Retreatment	16.05	321 (15.74)
Not available		39 (1.91)
	<i>total n = 2000</i>	
Geographical zone		
Central	1.10	22 (1.08)
Coastal	51.55	1029 (50.47)
Lake	17.94	358 (17.56)
Northern	20.19	403 (19.76)
S. Highlands	8.12	162 (7.95)
Western	0.50	10 (0.49)
Zanzibar	0.60	12 (0.59)
Not available		43 (2.11)
	<i>total n = 1996</i>	

149 IQR, interquartile range

150 Our study population consisted of TB patients ranging between the age of 2 and 82 years

151 with a median age of 28 years (interquartile range [IQR] 20–38). To further probe the age

152 distribution in the study population, we categorized the TB patients into five different age

153 groups (Table 1). We detected approximately three-quarters of the TB cases to occur among
154 the “young age” and “early adult” age groups. Our observation suggests that TB incidence in
155 Tanzania like in other high burden settings [30] is largely contributed by ongoing
156 transmission (rapid progression upon exposure to infection) as opposed to reactivation
157 (longer latency). Further, our findings corroborate with the national TB notification rates in
158 that about 10% of the TB cases are pediatric cases (< 15 years) [31].

159 Similar to other settings [1], we identified a higher proportion of male TB cases compared to
160 female TB cases. However, the male-to-female ratio in our study population is higher than
161 the national estimates for the two years of sampling (2.2:1 vs., 1.4:1). The striking gender
162 imbalance among TB cases seems to peak at adolescence onwards and is less pronounced
163 among pediatric TB cases (S1 Table). Additionally, a-third (32.2%, 517/1604) of the TB cases
164 with available HIV status were HIV co-infected. In contrast TB/HIV co-infected cases were
165 more likely to be female (44.5%, CI 38.3-50.7% vs., 25.8%, 95% CI 20.6-31.0%) which is
166 consistent with HIV being generally more prevalent in females than males [32]. We found
167 that our study population comprised 16.1% (321/2000) of TB retreatment cases, which was
168 four-fold higher than the overall countrywide notifications [31]. Finally, more than half
169 (51.6%, 1029/1996) of the TB patients in our study population were diagnosed in the Coastal
170 zone of Tanzania and about 40% were either diagnosed in the Lake and Northern zones. In
171 addition to higher TB notification rates, the three former mentioned geographical zones
172 contain the country’s zonal TB reference laboratories. The remaining 10% of the patients
173 were diagnosed in any of the remaining four geographical zones.

174 **Main MTBC lineages in Tanzania**

175 Using SNP-typing, we detected four of the seven known MTBC lineages (Fig 1), albeit at
176 varying proportions. In our study setting, Lineage 4 and Lineage 3 were the most frequent
177 (866, 42.5% and 681, 33.4%, respectively) followed by Lineage 1 (336, 16.5%). Lineage 2 was
178 the least frequent (92, 4.5%). The remaining 64 clinical isolates (3.1%) could not be assigned
179 into any of the MTBC lineages. Of the seven geographical zones, four (Coastal, Northern,
180 Lake and Southern Highlands) were highly represented with more than 100 clinical strains
181 each (Table 2). The distribution of the *M. tuberculosis* lineages varied within the
182 geographical zones (Fig 1 and S3 Fig). Our findings reveal that Lineage 1 strains were more
183 frequent in the Lake zone compared to the overall average frequency (20.9% vs. 16.8%),
184 whereas the frequency of Lineage 3 in this zone was lower (27.6% vs. 34.3%) compared to
185 other geographical zones. By contrast, Lineage 4 was the most predominant in all
186 geographical zones and showed relatively similar frequencies across the zones.

187

188 **Fig 1. MTBC lineages in Tanzania.** A. MTBC lineage classification of 2,039 nationwide clinical
189 strains. B. MTBC lineage frequencies and geographical distribution in Tanzania.

190

191 **Table 2.** *M. tuberculosis* lineage distribution across geographical regions in Tanzania

Geographical Zone	Lineage				Total
	L1 (%)	L2 (%)	L3 (%)	L4 (%)	
Central	8 (38.1)	2 (9.5)	4 (19)	7 (33.3)	21
Coastal	168 (16.8)	50 (5)	350 (35)	432 (43.2)	1000
Lake	72 (20.9)	12 (3.5)	95 (27.6)	165 (48)	344
Northern	52 (13.3)	22 (5.6)	145 (37)	173 (44.1)	392
S. Highlands	27 (16.9)	4 (2.5)	60 (37.5)	69 (43.1)	160
Western	0 (0)	1 (10)	4 (40)	5 (50)	10
Zanzibar	2 (18.2)	0 (0)	5 (45.5)	4 (36.4)	11
Total	329 (17)	91 (4.7)	663 (34.2)	855 (44.1)	1938

192 L1, Lineage 1; L2, Lineage 2; L3, Lineage 3; L4, Lineage 4

193 Sublineages

194 After we detected four main *M. tuberculosis* lineages, we next investigated the respective
195 sublineages within Lineage 1, 3 and 4 using spoligotyping. Lineage 2 strains were excluded
196 from this analysis since the global strains almost exclusively belong to one spoligotype
197 family, Beijing with almost identical fingerprint pattern. We identified 24 spoligotypes (SITs)
198 among the 107 clinical strains analyzed (S2 Table). Twenty six (24.3%) of the strains could
199 not be assigned to any of the existing spoligotypes in the SITVITWEB database and therefore
200 described as orphan spoligotypes. Several spoligotypes were identified within each of the
201 three lineages. Lineage 1 strains mainly belonged to EAI5 spoligotype. On the other hand,
202 CAS1_Kili was the most predominant spoligotype among the Lineage 3 strains. Within
203 Lineage 4 strains, LAM, T, and H families were detected and expectedly the LAM sublineage,
204 particularly LAM_ZWE was the most prevalent.

205 **Associations between lineages and patients' characteristics**

206 Having described the circulating main lineages of the *M. tuberculosis* we then assessed the
207 relationship between the lineages and patients' characteristics (Table 3). We detected a
208 higher proportion of female sex among TB patients infected with Lineage 2 (52.1%)
209 compared to those infected with the other three lineages (range from 31% to 34.5%, $p =$
210 0.009). Moreover, we observed that retreatment cases were frequently infected with
211 Lineage 2 strains (26.8%), which was twofold higher compared to Lineage 1 and 4 strains ($p <$
212 0.001). We found no evidence for association between lineages and patients' characteristics
213 such as age and HIV status (Table 3).

214 Lineage 2 has previously been associated with retreatment cases, drug resistance and lately
215 also with female sex [17,27]. We therefore investigated if similar associations exist in our
216 study population using a subset of TB cases with complete clinical and demographic
217 information ($n = 1515$). To assess these associations we performed logistic regression
218 analyses comparing Lineage 2 to all other lineages pooled together (Table 4). Our analyses
219 revealed Lineage 2 to be independently associated with female sex (adjusted odds ratio
220 [aOR] 2.25; 95% confidence interval [95% CI] 1.38 – 3.70, $p < 0.001$) and retreatment cases
221 (aOR 1.78; 95% CI 1.00 – 3.02, $p = 0.040$). We did not detect any association between the
222 lineages and patients' age and the HIV status.

223

224 **Table 3.** Frequency distribution of *M. tuberculosis* main lineages across patients'
 225 characteristic groups

Patient characteristics	Lineage			
	L1, n (%)	L2, n (%)	L3, n (%)	L4, n (%)
Age group				
Child age (< 15)	25 (9.3)	5 (7)	55 (10.6)	71 (10.8)
Young age (15 - 24)	76 (28.4)	21 (29.6)	153 (29.4)	202 (30.8)
Early adult (25 - 44)	124 (46.3)	37 (52.1)	265 (50.99)	294 (44.9)
Late adult (45 - 64)	36 (13.4)	6 (8.5)	39 (7.5)	80 (12.29)
Old age (> 65)	7 (2.6)	2 (2.8)	9 (1.7)	8 (1.2)
Sex				
Female	83 (31)	37 (52.1)	180 (34.5)	220 (33.6)
Male	185 (69)	34 (47.9)	341 (65.5)	435 (66.49)
HIV status				
Negative	181 (67.5)	45 (63.4)	341 (65.5)	452 (69)
Positive	87 (32.5)	26 (36.6)	180 (34.5)	203 (31)
Patient category				
New case	232 (86.6)	52 (73.2)	400 (76.8)	558 (85.2)
Retreatment	36 (13.4)	19 (26.8)	121 (23.2)	97 (14.8)
Total	268 (17.7)	71 (4.7)	521 (34.4)	655 (43.2)

226 L1, Lineage 1; L2, Lineage 2; L3, Lineage 3; L4, Lineage 4

227

228 **Table 4.** Associations of patients' clinical and demographic characteristics with *M.*
229 *tuberculosis* Lineage 2 (n = 71) compared to all other lineages (n = 1444)

Patient characteristics	Lineage 2	Unadjusted		Adjusted	
	n (%)	OR (95% CI)	p value	OR (95% CI)	p value
Age, median (IQR)	27 (20.5 – 38.5)			0.99 (0.98 – 1.01)	0.428
Female sex	37 (52.1)	2.17 (1.34–3.51)	0.002	2.25 (1.38 – 3.70)	<0.001
Retreatment case	19 (26.8)	1.71 (0.97–2.90)	0.052	1.78 (1.00 – 3.02)	0.040
HIV positive	26 (36.6)	0.98 (0.74 – 1.03)	0.915	1.03 (0.61 – 1.70)	0.91
Observations		1515		1515	

230 IQR, Interquartile range; OR, Odds ratio; 95% CI, 95% confidence interval.

231 **Mutations within *rpoB* gene in retreatment cases**

232 To investigate whether drug resistance was linked to a particular lineage, we included in
233 total 145 retreatment cases for drug resistance genotyping of the *rpoB* gene that confers
234 resistance to rifampicin. Out of these, 112 (77.2%) had no mutations compared to the H37Rv
235 reference gene and 16 (11%) contained at least one mutation, either synonymous (3/16) or
236 non-synonymous (13/16) (S4 Fig and S3 Table). We could not determine mutation status in
237 the *rpoB* gene of 17 (11.7%) retreatment cases due to PCR and sequencing failure. Among
238 the 13 strains detected with non-synonymous *rpoB* mutations, five belonged to Lineage 2,
239 four to Lineage 4, three to Lineage 3 and one was unclassified (S4 Table). Table 4 summarizes
240 the non-synonymous *rpoB* mutations detected.

241

242 **Table 4.** Detected mutations on the *rpoB* gene among the retreatment cases

rpoB mutation	Amino acid change	Comment
A1198G	T400A	reported
A1304T	D435V	reported
A1334T	H445L	reported
C1333G	H445D	reported
C1333T	H445Y	reported
T1289C	L430P	reported
A1442G	E481A	reported
C1294G	Q432E	reported
C1349T	S450L	reported

243

244

245 Discussion

246 In this study, we classified the countrywide collection of 2,039 *M. tuberculosis* isolates from
247 smear positive new and retreatment TB cases diagnosed between 2012 and 2013 in
248 Tanzania. Our findings show that the *M. tuberculosis* strains in Tanzania are diverse,
249 comprising four main phylogenetic lineages (Lineage 1–4) which occur throughout the
250 country. Specifically, we found that Lineage 4 was the most frequent nationwide, followed
251 by Lineage 3 and 1. Despite Lineage 2's recent global dissemination [15], it was the least
252 frequent in our study setting. Finally, our analysis on the relationship between *M.*
253 *tuberculosis* lineages and patients' characteristics revealed associations of Lineage 2 with
254 female sex and retreatment TB cases.

255 Among the 7 human–adapted MTBC lineages, Lineage 4 is the most broadly distributed and
256 occurs at high frequencies in Europe, the Americas and Africa [26,33]. We observe that TB
257 epidemics in Tanzania are also predominated by Lineage 4, which is regarded as the most
258 successful of MTBC lineages [33]. In general, the wide geographical range of Lineage 4 is
259 postulated to be driven by a combination of its enhanced virulence, high rates of human
260 migration linked to its spread and ultimately its ability to infect different human population
261 backgrounds [33,34]. In contrast, Lineage 1 and 3 are known to be mainly confined to the
262 rim of the Indian Ocean [9], which is consistent with our observation that nearly 50% of the
263 *M. tuberculosis* strains in Tanzania belong to these two lineages. This high prevalence of
264 Lineage 1 and 3 likely reflects the long-term migrations between Eastern Africa and the
265 Indian subcontinent [35]. In addition, the distribution and the frequency of Lineage 1 and 3
266 in the mainland away from the coast of Tanzania did not vary, suggesting spread via internal
267 migrations. Lineage 1 was proposed to have evolved in East Africa prior disseminating out of
268 the continent [12]. Based on this, one might expect higher frequencies of Lineage 1 in the

269 region. Instead, the so called “modern” (TbD1–) lineages (4 and 3 in this case) dominate in
270 Tanzania despite presumably being introduced into the African continent only post-contact
271 [33,36]. This illustrates the ability of “modern” lineages to thrive in co-existence with the
272 pre-existing “ancient” (TbD1+) lineages such as Lineage 1 in our case, perhaps because of the
273 comparably higher virulence [16,37]. The neighboring countries of Tanzania on the other
274 hand show comparable *M. tuberculosis* lineage composition [38,39], indicating common
275 demographic histories and ongoing exchanges that resulted into distinct *M. tuberculosis*
276 populations. The frequency of Lineage 2–Beijing in Tanzania, like in most parts of the
277 continent except for South Africa [38,39] is relatively low, despite the long-standing African-
278 Asian contacts [39]. Evidence from recent studies show that Lineage 2–Beijing was only
279 recently introduced into Africa [15,40].

280 The burden of TB disease is generally higher in males [1,41], rendering male sex as a
281 potential risk factor for TB. Furthermore, the male bias among TB patients is also observed in
282 settings with no obvious sex-based differences in health-seeking behavior [42]. Whilst we
283 show similar trends in this study setting overall, our findings reveal that the proportion of
284 females was higher among TB patients infected with Lineage 2. This finding is consistent with
285 several other previous studies conducted in different settings [17,27,43]. Social and
286 physiological factors predisposing males to higher risk of TB have been indicated [44]. On
287 the one hand, these include risk behaviors such as substance abuse (alcoholism, tobacco
288 smoking) and gender specific roles such as risk occupations (e.g., mining) that are male
289 dominated and known to increase the risk for TB. On the other hand, genetic makeup and
290 sex hormones might contribute to the differences in TB susceptibility among females and
291 males, as epidemiological and experimental studies have suggested female sex hormones to
292 be protective [44]. These observations would propose that the sex imbalance in TB to

293 emerge after the onset of puberty. Of note, we observe less sex imbalance in “child” age
294 group (<15 years) which also corroborates the national notification rates [31]. However, this
295 observation can be confounded by BCG vaccination which might be most effective in this age
296 group. Despite the high prevalence of HIV among young females in sub-Saharan Africa [32]
297 and HIV being the strongest risk factor for TB, TB burden remains higher in males. While
298 social and physiological aspects play an important role, findings from this study and others
299 previously conducted in Nepal and Vietnam [17,27] suggest that bacterial factors could
300 disrupt the trends towards male bias in TB, a finding which warrants further investigation.
301 Our hypothesis is that because of higher virulence, Lineage 2 strains are able to overcome
302 the resistance posed by female sex which could explain the less pronounced sex
303 imbalanced.

304 In addition to its association with female sex, Lineage 2 was also associated with retreatment
305 TB cases [43]. A retreatment case in our study population represented recurrent TB case
306 either due to relapse or reinfection. We hypothesized that this observation was possibly
307 linked to drug resistance, given the previous reported association between Lineage 2 and
308 drug resistance [45]. However, we detected only 9% (13/145) of strains among the
309 retreatment subset tested to contain mutations conferring resistance to rifampicin, five of
310 which belonged to Lineage 2. These findings would suggest that retreatment cases are
311 mainly driven by reinfection as opposed to treatment failure or relapse.

312 Finally, as evidenced by the age distribution of TB cases in our study setting, recent or
313 ongoing transmission in high burden countries is the main contributor to the TB burden
314 rather than disease reactivation [30]. Additionally, an association with young age has been
315 employed as an epidemiological proxy for highly transmissible strains and faster rates of
316 disease progression [46,47]. In this study, we did not detect any differences in median age

317 of TB patients infected with different lineages (S5 Fig), an observation that could speak for
318 high ongoing transmission rates in general, irrespective of lineage.

319 Our study is limited by focusing on a convenient collection of archived *M. tuberculosis*
320 clinical isolates (representing 1.6% of all TB cases in 2012 and 2013) sampled from TB cases
321 as part of countrywide drug resistance surveillance. Therefore, the strength or lack of
322 associations between lineages and patients' characteristics could likely be affected by the
323 sampling. In addition, most of the geographical zones were underrepresented which could in
324 turn underestimate the respective regional lineage composition and the overall countrywide
325 distribution. Systematic sampling would allow for better resolution on the distribution
326 patterns, the frequencies and on epidemiological consequences of *M. tuberculosis* lineages,
327 which might partially determine the regional specific epidemics.

328 In conclusion, this study addresses the countrywide *M. tuberculosis* population structure
329 based on SNP-typing. We show that *M. tuberculosis* population in Tanzania is diverse with
330 four of the seven known lineages detected. This study sets the stage for further in depth
331 investigations on epidemiological impact of *M. tuberculosis* lineages in Tanzania.

332

333 **Acknowledgements**

334 We would like to thank the National Tuberculosis Leprosy Programme (NTLP) through the
335 Central Reference Laboratory (CTRL) for permission to use the *M. tuberculosis* isolate
336 collection for this study.

337

338 References

- 339 1. WHO. Global tuberculosis report. Geneva: World Health Organization. 2017.
- 340 2. National Tuberculosis and Leprosy Programme (NTLP). Annual report for 2016. Dar es
341 Salaam. 2016;1–48.
- 342 3. Nagu TJ, Aboud S, Mwiru R, Matee M, Fawzi W, Mugusi F. Multi Drug and Other Forms
343 of Drug Resistant Tuberculosis Are Uncommon among Treatment Naïve Tuberculosis
344 Patients in Tanzania. Surolia A, editor. PLoS One. 2015 Apr 7;10(4):e0118601.
- 345 4. Ministry of Health and Social Welfare . Dar es Salaam. First tuberculosis prevalence
346 survey in the United Republic of Tanzania. 2013;
- 347 5. Rieder HL. Epidemiologic Basis of Tuberculosis Control First edition 1999. 1999.
- 348 6. Houben RMGJ, Dodd PJ. The Global Burden of Latent Tuberculosis Infection: A Re-
349 estimation Using Mathematical Modelling. Metcalfe JZ, editor. PLOS Med. 2016 Oct
350 25;13(10):e1002152.
- 351 7. Koul A, Arnoult E, Lounis N, Guillemont J, Andries K. The challenge of new drug
352 discovery for tuberculosis. Nature. 2011 Jan 27;469(7331):483–90.
- 353 8. Comas I, Gagneux S. The Past and Future of Tuberculosis Research. Manchester M,
354 editor. PLoS Pathog. 2009;5(10):e1000600.
- 355 9. Gagneux S. Ecology and evolution of Mycobacterium tuberculosis. Nat Rev Microbiol.
356 2018 Feb 19;16(4):202–13.
- 357 10. Gagneux S, DeRiemer K, Van T, Kato-Maeda M, de Jong BC, Narayanan S, et al.
358 Variable host-pathogen compatibility in Mycobacterium tuberculosis. Proc Natl Acad
359 Sci. 2006 Feb 21;103(8):2869–73.
- 360 11. Fenner L, Egger M, Bodmer T, Furrer H, Ballif M, Battegay M, et al. HIV Infection
361 Disrupts the Sympatric Host–Pathogen Relationship in Human Tuberculosis. Gibson G,
362 editor. PLoS Genet. 2013 Mar 7;9(3):e1003318.
- 363 12. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, et al. Out-of-Africa
364 migration and Neolithic coexpansion of Mycobacterium tuberculosis with modern
365 humans. Nat Genet. 2013 Oct 1;45(10):1176–82.
- 366 13. Coscolla M. Biological and Epidemiological Consequences of MTBC Diversity. In: Strain
367 variation in the Mycobacterium tuberculosis complex:Its role in biology, epidemiology
368 and control. Springer, Cham; 2017. p. 95–116.
- 369 14. Hanekom M, van der Spuy GD, Streicher E, Ndabambi SL, McEvoy CRE, Kidd M, et al. A
370 Recently Evolved Sublineage of the Mycobacterium tuberculosis Beijing Strain Family
371 Is Associated with an Increased Ability to Spread and Cause Disease. J Clin Microbiol.
372 2007 May 1;45(5):1483–90.
- 373 15. Cowley D, Govender D, February B, Wolfe M, Steyn L, Evans J, et al. Recent and Rapid
374 Emergence of W-Beijing Strains of Mycobacterium tuberculosis in Cape Town, South

- 375 Africa. *Clin Infect Dis*. 2008 Nov 15;47(10):1252–9.
- 376 16. Stavrum R, PrayGod G, Range N, Faurholt-Jepsen D, Jeremiah K, Faurholt-Jepsen M, et al. Increased level of acute phase reactants in patients infected with modern
377 Mycobacterium tuberculosis genotypes in Mwanza, Tanzania. *BMC Infect Dis*. 2014
378 Dec 5;14(1):309.
- 380 17. Holt KE, McAdam P, Thai PVK, Thuong NTT, Ha DTM, Lan NN, et al. Frequent
381 transmission of the Mycobacterium tuberculosis Beijing lineage and positive selection
382 for the EsxW Beijing variant in Vietnam. *Nat Genet*. 2018 Jun 21;50(6):849–56.
- 383 18. Eldholm V, Matee M, Mfinanga SGM, Heun M, Dahle UR. A first insight into the
384 genetic diversity of Mycobacterium tuberculosis in Dar es Salaam, Tanzania, assessed
385 by spoligotyping. *BMC Microbiol*. 2006 Sep 13;6(1):76.
- 386 19. Kibiki GS, Mulder B, Dolmans WM, de Beer JL, Boeree M, Sam N, et al. M. tuberculosis
387 genotypic diversity and drug susceptibility pattern in HIV- infected and non-HIV-
388 infected patients in northern Tanzania. *BMC Microbiol*. 2007 May 31;7(1):51.
- 389 20. Mfinanga SGM, Warren RM, Kazwala R, Kahwa A, Kazimoto T, Kimaro G, et al. Genetic
390 profile of Mycobacterium tuberculosis and treatment outcomes in human pulmonary
391 tuberculosis in Tanzania. *Tanzan J Health Res*. 2014 Apr;16(2):58–69.
- 392 21. Mbugi E V., Katale BZ, Siame KK, Keyyu JD, Kendall SL, Dockrell HM, et al. Genetic
393 diversity of Mycobacterium tuberculosis isolated from tuberculosis patients in the
394 Serengeti ecosystem in Tanzania. *Tuberculosis*. 2015 Mar;95(2):170–8.
- 395 22. Comas I, Homolka S, Niemann S, Gagneux S. Genotyping of Genetically Monomorphic
396 Bacteria: DNA Sequencing in Mycobacterium tuberculosis Highlights the Limitations of
397 Current Methodologies. Litvintseva AP, editor. *PLoS One*. 2009 Nov 12;4(11):e7815.
- 398 23. Fenner L, Malla B, Ninet B, Dubuis O, Stucki D, Borrell S, et al. “Pseudo-Beijing”:
399 Evidence for Convergent Evolution in the Direct Repeat Region of Mycobacterium
400 tuberculosis. Sechi LA, editor. *PLoS One*. 2011 Sep 13;6(9):e24737.
- 401 24. Stucki D, Malla B, Hostettler S, Huna T, Feldmann J, Yeboah-Manu D, et al. Two New
402 Rapid SNP-Typing Methods for Classifying Mycobacterium tuberculosis Complex into
403 the Main Phylogenetic Lineages. Mokrousov I, editor. *PLoS One*. 2012 Jul
404 20;7(7):e41253.
- 405 25. Kamerbeek J, Schouls L, Kolk A, van Agterveld M, van Soolingen D, Kuijper S, et al.
406 Simultaneous detection and strain differentiation of Mycobacterium tuberculosis for
407 diagnosis and epidemiology. *J Clin Microbiol*. 1997 Apr;35(4):907–14.
- 408 26. Demay C, Liens B, Burguière T, Hill V, Couvin D, Millet J, et al. SITVITWEB – A publicly
409 available international multimarker database for studying Mycobacterium
410 tuberculosis genetic diversity and molecular epidemiology. *Infect Genet Evol*. 2012
411 Jun;12(4):755–66.
- 412 27. Malla B, Stucki D, Borrell S, Feldmann J, Maharjan B, Shrestha B, et al. First Insights
413 into the Phylogenetic Diversity of Mycobacterium tuberculosis in Nepal. Sola C, editor.
414 *PLoS One*. 2012 Dec 26;7(12):e52297.

- 415 28. Staden R. The Staden sequence analysis package. *Mol Biotechnol.* 1996 Jun;5(3):233–
416 41.
- 417 29. R Core Team. R: A Language and Environment for Statistical Computing. 2018.
- 418 30. Yates TA, Khan PY, Knight GM, Taylor JG, McHugh TD, Lipman M, et al. The
419 transmission of Mycobacterium tuberculosis in high burden settings. *Lancet Infect Dis.*
420 2016;16(2):227–38.
- 421 31. National Tuberculosis and Leprosy Programme (NTLP). Annual report for 2013. Dar es
422 Salaam. 2013.
- 423 32. Hegdahl HK, Fylkesnes KM, Sandøy IF. Sex Differences in HIV Prevalence Persist over
424 Time: Evidence from 18 Countries in Sub-Saharan Africa. Faragher EB, editor. *PLoS*
425 *One.* 2016 Feb 3;11(2):e0148502.
- 426 33. Stucki D, Brites D, Jeljeli L, Coscolla M, Liu Q, Trauner A, et al. Mycobacterium
427 tuberculosis lineage 4 comprises globally distributed and geographically restricted
428 sublineages. *Nat Genet.* 2016 Dec 31;48(12):1535–43.
- 429 34. Coscolla M, Gagneux S. Consequences of genomic diversity in Mycobacterium
430 tuberculosis. *Semin Immunol.* 2014 Dec;26(6):431–44.
- 431 35. O’Neill MB, Shockey AC, Zarley A, Aylward W, Eldholm V, Kitchen A, et al. Lineage
432 specific histories of Mycobacterium tuberculosis dispersal in Africa and Eurasia.
433 *bioRxiv.* 2018 Aug 1;210161.
- 434 36. Comas I, Hailu E, Kiros T, Bekele S, Mekonnen W, Gumi B, et al. Population Genomics
435 of Mycobacterium tuberculosis in Ethiopia Contradicts the Virgin Soil Hypothesis for
436 Human Tuberculosis in Sub-Saharan Africa. *Curr Biol.* 2015 Dec 21;25(24):3260–6.
- 437 37. Portevin D, Gagneux S, Comas I, Young D. Human Macrophage Responses to Clinical
438 Isolates from the Mycobacterium tuberculosis Complex Discriminate between Ancient
439 and Modern Lineages. Bessen DE, editor. *PLoS Pathog.* 2011 Mar 3;7(3):e1001307.
- 440 38. Mbugi E V., Katale BZ, Streicher EM, Keyyu JD, Kendall SL, Dockrell HM, et al. Mapping
441 of Mycobacterium tuberculosis Complex Genetic Diversity Profiles in Tanzania and
442 Other African Countries. Sreevatsan S, editor. *PLoS One.* 2016 May 5;11(5):e0154571.
- 443 39. Chihota VN, Niehaus A, Streicher EM, Wang X, Sampson SL, Mason P, et al. Geospatial
444 distribution of Mycobacterium tuberculosis genotypes in Africa. Arez AP, editor. *PLoS*
445 *One.* 2018 Aug 1;13(8):e0200632.
- 446 40. Rutaiwa LK, Menardo F, Stucki D, Gygli SM, Ley SD, Malla B, et al. Multiple
447 introductions of the Mycobacterium tuberculosis Lineage 2 Beijing into Africa over
448 centuries. *bioRxiv.* 2018 Sep 10;413039.
- 449 41. Guerra-Silveira F, Abad-Franch F. Sex Bias in Infectious Disease Epidemiology: Patterns
450 and Processes. Nishiura H, editor. *PLoS One.* 2013 Apr 24;8(4):e62390.
- 451 42. Rhines AS. The role of sex differences in the prevalence and transmission of
452 tuberculosis. *Tuberculosis.* 2013 Jan 1;93(1):104–7.

- 453 43. Buu TN, Huyen MN, Lan NTN, Quy HT, Hen N V, Zignol M, et al. The Beijing genotype is
454 associated with young age and multidrug-resistant tuberculosis in rural Vietnam. *Int J*
455 *Tuberc Lung Dis.* 2009 Jul;13(7):900–6.
- 456 44. Nhamoyebonde S, Leslie A. Biological Differences Between the Sexes and
457 Susceptibility to Tuberculosis. *J Infect Dis.* 2014 Jul 15;209(suppl 3):S100–6.
- 458 45. Borrell S, Gagneux S. Infectiousness, reproductive fitness and evolution of drug-
459 resistant *Mycobacterium tuberculosis* [State of the art]. *Int J Tuberc Lung Dis.*
460 2009;13(12):1456–1466.
- 461 46. de Jong BC, Hill PC, Aiken A, Awine T, Antonio M, Adetifa IM, et al. Progression to
462 Active Tuberculosis, but Not Transmission, Varies by *Mycobacterium tuberculosis*
463 Lineage in The Gambia. *J Infect Dis.* 2008 Oct 1;198(7):1037–43.
- 464 47. Borgdorff MW, van Soolingen D. The re-emergence of tuberculosis: what have we
465 learnt from molecular epidemiology? *Clin Microbiol Infect.* 2013 Oct;19(10):889–901.
- 466

467 **Supporting information**

468 **S1 Fig. Study population flowchart.**

469 **S2 Fig. Patients' data included in the study.** Proportion of observed and missing data for the
470 variables included in the study

471 **S3 Fig. MTBC lineage proportions.** Distribution of *M. tuberculosis* lineages across different
472 regions of Tanzania. Size of the circle is proportional to the number of isolates analyzed from
473 the regions.

474 **S4 Fig. Flowchart of genotyped strains for *rpoB* mutations.** A subset of *M. tuberculosis*
475 strains from retreatment cases included for *rpoB* drug resistance genotyping

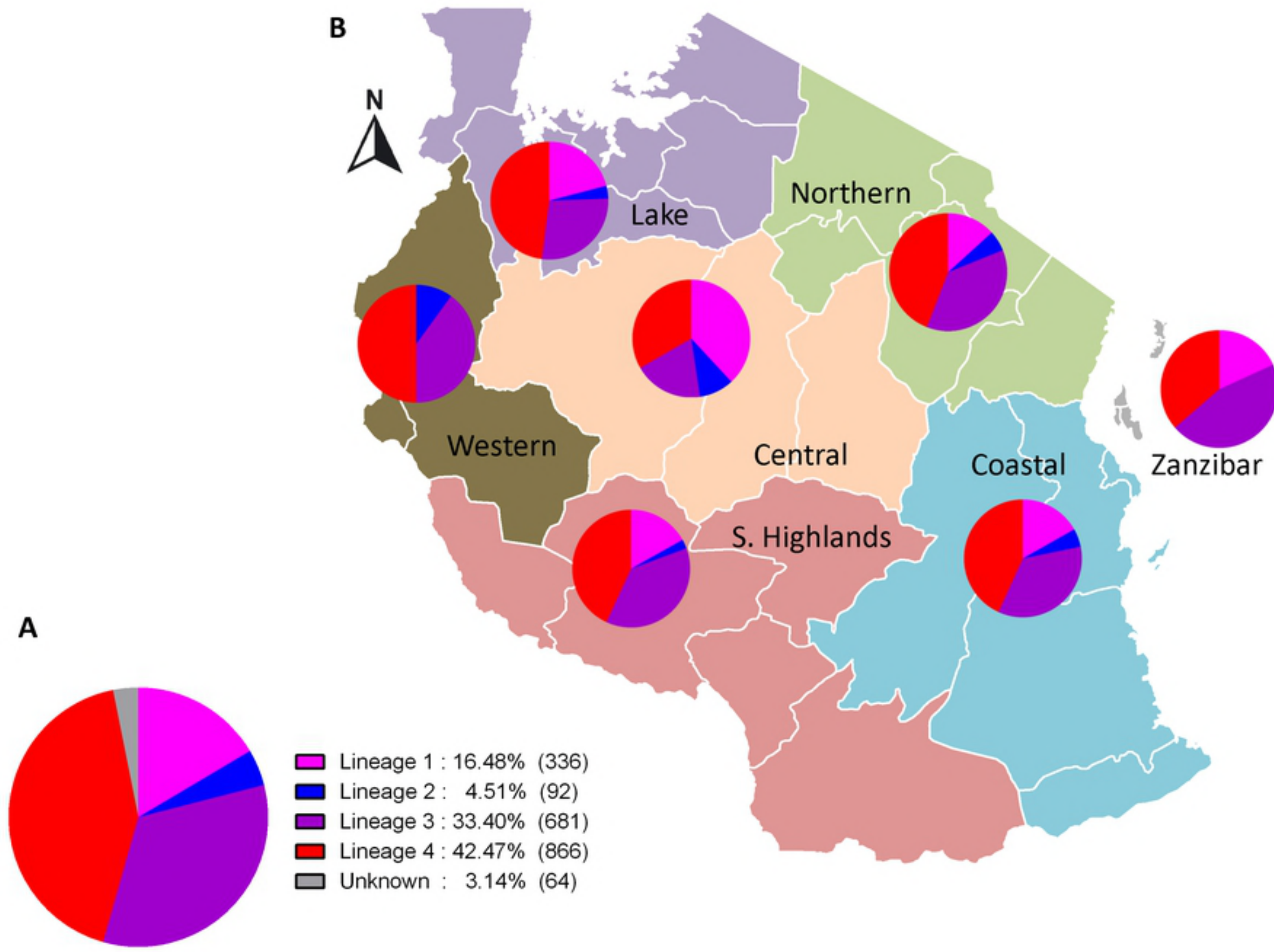
476 **S5 Fig. Patients' age distribution across MTBC lineages.** The age distributions of TB patients
477 grouped by infecting MTBC lineage.

478 **S1 Table. Sex distribution across different age groups of TB patients.**

479 **S2 Table. Spoligotype patterns of a subset of *M. tuberculosis* clinical strains**

480 **S3 Table. Mutations detected in the *rpoB* gene**

481 **S4 Table. Distribution of *rpoB* mutations across the four MTBC lineages**



Figure