1   **Using genomic epidemiology and geographic activity spaces to investigate tuberculosis**

2   **outbreaks in Botswana**

3

4   **Chelsea R. Baker, Ivan Barilar, Leonardo S. de Araujo, Daniel M. Parker, Kimberly**

5   **Fornace, Patrick K. Moonan, John E. Oeltmann, James L. Tobias, Volodymyr M. Minin,**

6   **Chawangwa Modongo, Nicola M. Zetola[1], Stefan Niemann[1], and Sanghyuk S. Shin[1]**

7   Author affiliations: University of California, Irvine, California, USA (C.R. Baker, D.M. Parker,

8   V.M. Minin, S.S. Shin); Forschungszentrum, Borstel, Germany (I. Barilar, L.S. de Araujo, S.

9   Niemann); National University of Singapore (K. Fornace); US Centers for Disease Control and

10  Prevention, Atlanta, Georgia, USA (P.K. Moonan, J.E. Oeltmann, J.L. Tobias,);  Botswana–

11  Upenn Partnership, Gaborone, Botswana/ Victus Global Botswana Organisation, Gaborone,

12  Botswana (C. Modongo, N.M. Zetola)

13

14  [1]These senior authors contributed equally to this article.

15  **Corresponding author:**

16  Sanghyuk Shin, PhD

17  Associate Professor, Sue & Bill Gross School of Nursing

18  Email: ssshin2@uci.edu

19  Phone: 949-576-8675

20

21  **Key words:** Tuberculosis transmission, spatial analysis, activity space, whole genome

22  sequencing, geographic heterogeneity, outbreak, infectious disease control

23    **Abstract**

24    <u>Background</u>

25    The integration of genomic and geospatial data into infectious disease transmission analyses

26    typically includes residential locations and excludes other activity spaces where transmission

27    may occur (*e.g.* work, school, or social venues). The objective of this analysis was to explore

28    residential as well as other activity spaces of tuberculosis (TB) outbreaks to identify potential

29    geospatial 'hotspots' of transmission.

30    <u>Methods</u>

31    We analyzed data that included geospatial coordinates for residence and other activity spaces

32    collected during 2012–2016 for the Kopanyo Study, a population-based study of TB transmission

33    in Botswana. We included participants with results from whole genome sequencing conducted on

34    archived samples from the original study. We used a spatial log-Gaussian Cox process model to

35    detect core areas of increased activity spaces of individuals belonging to TB outbreaks

36    (genotypic groups with ≤5 single-nucleotide polymorphisms), which we compared to ungrouped

37    participants (those not in a genotypic group of any size).

38    <u>Findings</u>

39    We analyzed data collected from 636 participants, including 70 participants belonging to six

40    outbreak groups with a combined total of 293 locations, and 566 ungrouped participants with a

41    combined total of 2289 locations. Core areas of activity space for each outbreak group were

42    geographically distinct, and we found evidence of localized transmission in four of six outbreaks.

43    For most of the outbreaks, including activity space data led to the detection of larger areas of

44    higher spatial intensity and more focal points compared to residential location alone.

45    <u>Interpretation</u>

46    Geospatial analysis using activity space data (social gathering places as well as residence) may

47    lead to improved understanding of areas of infectious disease transmission compared to using

48    residential data alone.

49    <u>Funding</u>

52

## Background

54    Tuberculosis (TB) remains among the leading causes of death due to infectious illness, despite

55    being a preventable and curable disease[1]. In 2022, over 10 million people became sick with TB,

56    and 1.3 million died[1]. Progress toward TB elimination has been slow and many targets set by the

57    World Health Organization (WHO) have not been met[1]. New strategies and tools are needed for

58    TB prevention[1]. In high-burden settings, where a substantial portion of disease incidence is due

59    to recent infection, interventions to stop ongoing transmission are especially important[1–5].

60

61    A promising tool is the integration of geospatial and pathogen genomic data. Pathogen whole

62    genome sequencing (WGS) can be used to identify closely related *M. tuberculosis* isolates and

63    help reconstruct likely transmission chains[6]. Geographic and genomic data can be combined to

64    help detect areas of sustained transmission, locate outbreaks, and investigate the geographic

65    range of different strains[7,8]. Spatial analysis of WGS data can help identify high-risk areas that

66    could be targeted for public health interventions to interrupt ongoing transmission[3–7,9–15].

67    Geographically targeted interventions have shown promise as an effective and cost-efficient

68    strategy for reducing TB incidence in high-burden, low-resource settings[16–19].

69

70    However, an important limitation to many studies employing this strategy is geospatial analysis

71    based solely on residential location, which excludes locations in the community where

72    transmission may occur[12,20–22]. An alternative approach is to analyze "activity space," which

73    includes the places one routinely occupies during day to day life[23–25]. For example, this may

74    include residential as well as community sites such as workplaces, markets, places of worship, or

75    other social gathering places[23,24]. This approach has the potential to lead to more accurate

76    detection of high-risk areas compared to analysis of residential locations alone[24,26].

77

78    We previously conducted a descriptive study of geospatial residential data and WGS data from a

79    population-based study of TB transmission in Botswana found evidence that TB outbreaks

80    displayed distinct geographic characteristics[27]. The objectives of the current analysis were to use

81    spatial statistical modeling to 1) identify geographic characteristics of the collective activity

82    space (residential as well as social gathering locations) of each outbreak group, and 2) identify

83    potential 'hotspot' areas of activity space associated with each outbreak, which may represent

84    areas of increased risk for transmission.

85

86    **Methods**

87    <u>Study design and setting</u>

88    We analyzed data collected during 2012–2016 for the Kopanyo Study, a population-based study

89    of TB transmission in Botswana, a country in southern Africa with a high burden of TB and

90    TB/HIV co-infection[1,5,28]. Participants were recruited at multiple local health clinics in two

91    districts: Gaborone, the urban center and capital city, and Ghanzi, a rural district several hundred

92    kilometers away[5,28]. During the five years before the study, TB incidence was 440–470

93    cases/100,000 persons in Gaborone, which had a total population of 354,380, and 722

94    cases/100,000 persons in Ghanzi, which had a population of 44,100 (12,179 in Ghanzi town)[5,28].

95

96    Study participants included men and women of all ages with TB disease who were sequentially

97    enrolled by date of diagnosis[5,28]. Those who had already received TB treatment for >14 days,

98    prisoners, and patients who declined to participate were excluded[5,28]. At least 1 sputum sample

99    was collected from each participant for bacterial culture[5,28]. Clinical and demographic data were

100   collected through in-person interviews and medical record review[5,28].

101

102   Data gathered during participant intake interviews included high resolution geospatial data for

103   activity space, which included home residence and social gathering places[5,28]. Participants were

104   asked about residential location as well as social gathering places (e.g. workplaces, schools,

105   markets, places of worship, alcohol venues etc.) frequented during their potential infectious

106   period (up to 12 months prior to treatment initiation)[5,28]. Geographic coordinates (latitude and

107   longitude recorded using the WGS 84 projection system with 1.1-m precision) for locations were

108   obtained using global positioning system (GPS) devices during site visits, or by geocoding

109   addresses using a reference layer created by manually relocating addresses in satellite imagery

110   using Google Maps, OpenStreetMap, and ArcGIS[5,28].

111

112   WGS

113   Whole genome sequencing was conducted on DNA samples archived from the original study

114   with sufficient quantities of DNA (>0.05 ng/μL) for analysis. Closely related *M. tuberculosis*

115   isolates were identified bioinformatically using a single linkage clustering algorithm. We

116   considered clusters of isolates with ≤5 single-nucleotide polymorphisms (SNPs) to indicate

117   recent transmission and clusters of ≥10 persons to be outbreaks. Further details of this procedure

118   are outlined in a separate analysis[27].

119

120   Spatial modeling of activity space

121 Participants eligible for the current analysis included those with WGS data, GPS coordinates,

122 and sociodemographic data for age, sex, income, and HIV status available. We focused our

123 current analysis on outbreak groups that had at least 10 activity space locations (collectively

124 among all their participants) within greater Gaborone, an area of approximately 27 km x 24 km

125 that includes the capital city and its surrounding suburbs. We also included genotypically

126 ungrouped participants as a comparison group. For model fitting purposes, a very small jitter was

127 introduced to location coordinates to avoid duplicate points (roughly on the scale of different

128 areas of the same building, ranging from approximately <1 to 10 meters).

129

130 We conducted a preliminary analysis to compare the geographic distribution of participants with

131 WGS data available to the total study population from the Kopanyo Study to rule out geographic

132 sampling bias. We estimated the geographic median center (a centralized point that minimizes

133 the distance to all other points), and directional distribution (which calculates the standard

134 deviation of points along both the X and Y axes) for both groups of participants and found nearly

135 identical results, indicating that participants with WGS data were geographically representative

136 of the larger study population. This analysis was performed using ArcGIS[29].

137

138 <u>Model description</u>

139 We used a spatial log-Gaussian Cox process (LGCP) to model the spatial intensity (average

140 number of points per unit area) of activity spaces of participants belonging to each outbreak

141 group ('cases') and of genotypically ungrouped participants (those not in an identified genotypic

142 group of any size, 'controls'). LGCPs are a flexible class of models for spatial point processes

143 where spatial intensity may vary across the study region[30,31]. A spatial random effect can be

144  incorporated to account for spatial correlation in the data and identify spatial patterns not

145  explained by other variables[33,34,37,40]. This technique offers a model-based approach for

146  estimating utilization distributions, which are probability density functions that can be mapped to

147  highlight areas with increased geographic concentrations of points (e.g. activity space locations)

148  to help characterize use of space[26,32]. To adapt the modeling framework to an activity space

149  context where each individual may be associated with multiple point locations, observations can

150  be treated as cumulative 'encounters' over specified time periods[32]. We considered each point to

151  represent an 'encounter' in space corresponding to a potential TB exposure, and estimated

152  intensity surfaces for cumulative exposures over the entire study period.

153

154  LGCPs fit well in a Bayesian hierarchical modeling framework, and various tools can be used for

155  this approach[33–35]. We used integrated nested Laplace approximation (INLA), a flexible and

156  computationally efficient method for approximate Bayesian inference for latent Gaussian

157  models, which include LGCPs[33,34,36–38]. We implemented this using the R-INLA package[39]. We

158  modeled the spatial random effect as a Gaussian random field (GRF) with Matérn

159  covariance[37,38,40]. We used the stochastic partial differential equation (SPDE) approach in R-

160  INLA to approximate the GRF[37,38,40]. We specified the SPDE model using penalized complexity

161  priors that were vaguely informative about the underlying spatial process (prior probability of

162  0.05 that the ranges of the fields were less than 0.5 km and prior probability of 0.05 that the

163  standard deviation was greater than 10).

164

165  Under the LGCP framework, we used a joint modeling approach to incorporate a shared spatial

166  term (obtained by jointly estimating the intensity of both cases and controls), as well as a unique

167    spatial term estimated for cases in each group[38,41]. Using this approach, posterior mean estimates

168    of the spatial random effect for cases represent variation in intensity not accounted for by the

169    spatial distribution of controls[41]. We did this to help identify areas with relatively high

170    concentrations of activity spaces associated with individual outbreak groups, while attempting to

171    account for baseline use of space (as some locations tend to be frequented by people more often

172    in general). Areas with a high density of activity spaces frequented by people belonging to the

173    same outbreak group could potentially represent areas associated with an increased risk of recent

174    transmission. We fit a version of the model that included just the shared spatial term (model 0),

175    and a version of the model that included the shared spatial term as well as unique spatial terms

176    estimated for each outbreak group individually (model 1).  We also conducted a sensitivity

177    analysis using subsets of the data with 70 and 140 randomly selected ungrouped participants as

178    controls to examine whether spatial patterns were sensitive to size of the control group.

179

180    We projected posterior mean estimates of the spatial effect (i.e. the effect of spatial location on

181    the intensity of activity spaces, represented by the spatial random field) for each outbreak group

182    onto maps of the study area in order to visualize how it varied across the region, and to identify

183    areas of increased or decreased (different than zero) values not explained by the spatial

184    distribution of controls[41]. Estimated values (displayed on the internal linear predictor scale)

185    represent the contribution of the spatial random effect to the response (spatial intensity) after

186    accounting for other fixed and random effects in the model. In addition, we reported posterior

187    mean estimates for the range (distance at which spatial correlation falls close to zero) and

188    variance of the spatial effect for each outbreak group[36].

189

190    We also projected posterior mean estimates for predicted spatial intensity values (fitted values of

191    the response at prediction locations, obtained by exponentiating the linear predictor), in order to

192    visualize patterns of spatial intensity of activity spaces for each outbreak group[36,38].

193

194    In addition, we calculated exceedance probabilities and projected these onto maps of the study

195    area to identify areas where estimated spatial effect for each outbreak group had a high

196    probability (0.95) of being greater than zero, representing high-confidence areas where the

197    spatial effect for cases was above the baseline that could be accounted for by the spatial

198    distribution of controls[41]. We also calculated exceedance probabilities to identify high-

199    confidence areas where the estimated spatial intensity was in the top ten percent of estimated

200    mean values for each group, representing 'core areas' or 'hotspots' of that group's collective

201    activity space[32].

202

203    We also generated exceedance probability maps based only on location of participant residence

204    (using the same threshold values as the full analysis) to compare high-risk areas identified using

205    activity space analysis vs. home location alone. We projected these exceedance probabilities onto

206    interactive maps of the study area for each outbreak group.

207

208    Map visualization was performed using the R packages raster, terra, sf, ggplot2, ggspatial, and

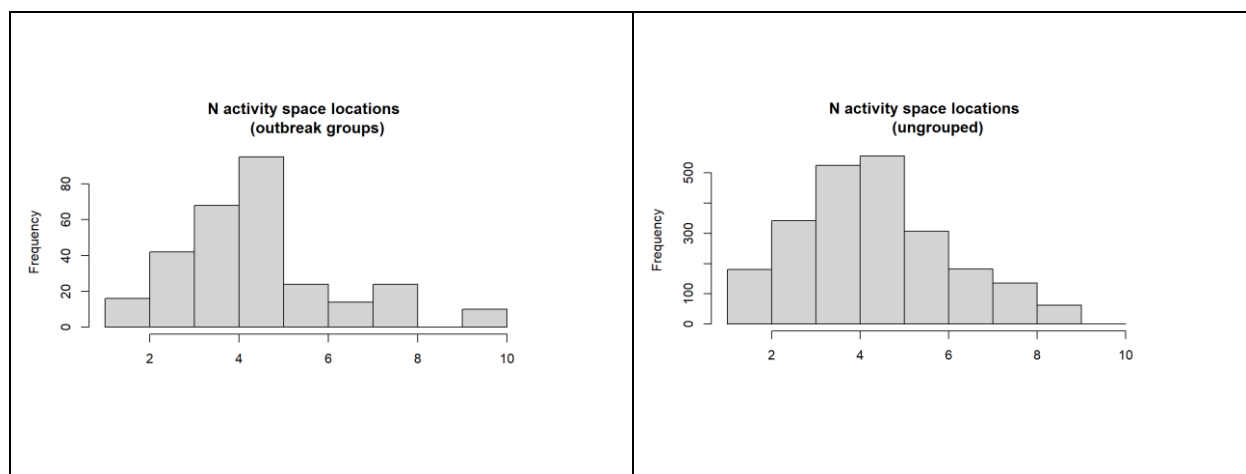209    leaflet.

210

211    **Results**

212 Participants

213 A total of 1426 participants had WGS data available, of which 1425 had GPS coordinates

214 available for at least one activity space location (home or social gathering place). Participants

215 with and without WGS data had similar sociodemographic characteristics in terms of age, sex,

216 HIV status, and income. Eight genotypic groups had 10 participants or more and were considered

217 outbreaks. Six out of eight outbreaks (genotypic groups with 10 participants or more) had at least

218 10 activity spaces (collectively among all their participants) in greater Gaborone.

219 A total of 636 participants with activity spaces in greater Gaborone met criteria for the current

220 analysis, including 70 participants belonging to six outbreak groups with a combined total of 293

221 locations, and 566 ungrouped participants with a combined total of 2289 locations.

222

223 Each participant had between one and 10 activity space locations, and the median number of

224 locations (n=4) was the same for both grouped and ungrouped participants (Figure 1). Median

225 number of activity space locations was the same (n=4) by gender and HIV status, though was

226 slightly lower for participants with no income (n=3) than participants with any income (n=4),

227 which could reflect an increased number of activity spaces among participants who were

228 employed. Among participants with more than one location, the maximum distance between any

229 two of their activity spaces ranged from <0.5 km to 21.2 km (median 6.2 km) for ungrouped

230 participants and <0.5 km to 21.7 km (median 4.3 km) for participants in outbreak groups

231 (supplementary figure 1).

232

233 Figure 1. Histograms for distribution of number of activity space locations per participants for

234 outbreak and genotypically ungrouped participants.

236 Among genotypically ungrouped participants, the median age was 35 years (IQR: 28–42), just

237 over half were male, about one quarter reported no income, and nearly 65% were diagnosed with

238 TB-HIV coinfection (Table 1). Among participants in the six genotypic groups, median age

239 ranged from 30 years (Group A) to 39 years (Group G) (Table 1). Participants in Group G were

240 exclusively male, while Group C alone was majority female (75%). Group D had the highest

241 proportion of participants diagnosed with TB-HIV coinfection (9 of 11; 91%). The percentage of

242 participants reporting no income ranged from 18% in Group D to 58% in Groups C and E (Table

243 1).

244

245 Table 1. Characteristics of study participants (N = 636) by outbreak group (genotypic group ≤ 5

246 SNP), Gaborone, Botswana, 2012-2016

| | A (N=22) | C (N=12) | D (N=11) | E (N=12) | G (N=9) | H (N=4) | Ungrouped (N=566) |
|---|---|---|---|---|---|---|---|
| Total locations | 81 | 45 | 53 | 54 | 44 | 16 | 2289 |
| Gender | | | | | | | |

| | A (N=22) | C (N=12) | D (N=11) | E (N=12) | G (N=9) | H (N=4) | Ungrouped (N=566) |
|---|---|---|---|---|---|---|---|
| Female | 11 (50.0%) | 9 (75.0%) | 5 (45.5%) | 3 (25.0%) | 0 (0%) | 2 (50.0%) | 264 (46.6%) |
| Male | 11 (50.0%) | 3 (25.0%) | 6 (54.5%) | 9 (75.0%) | 9 (100%) | 2 (50.0%) | 302 (53.4%) |
| **Age** | | | | | | | |
| Median [Q1,Q3] | 29 [24, 37] | 31 [29, 36] | 33 [31, 42] | 35 [29, 40] | 39 [35, 42] | 24 [20, 38] | 35 [28, 42] |
| **HIV Status** | | | | | | | |
| Neg | 10 (45.5%) | 5 (41.7%) | 1 (9.1%) | 6 (50.0%) | 4 (44.4%) | 3 (75.0%) | 203 (35.9%) |
| Pos | 12 (54.5%) | 7 (58.3%) | 10 (90.9%) | 6 (50.0%) | 5 (55.6%) | 1 (25.0%) | 363 (64.1%) |
| **Income** | | | | | | | |
| Any | 16 (72.7%) | 5 (41.7%) | 9 (81.8%) | 5 (41.7%) | 7 (77.8%) | 2 (50.0%) | 417 (73.7%) |
| None | 6 (27.3%) | 7 (58.3%) | 2 (18.2%) | 7 (58.3%) | 2 (22.2%) | 2 (50.0%) | 149 (26.3%) |

247

## Estimated spatial effects

249 Model 1 (shared and group-specific spatial terms) had a lower DIC (-12956.58) than model 0

250 (shared spatial terms only, DIC -12853.84), supporting the presence of spatial variation among

251 genotypic groups not accounted for by the spatial distribution of activity spaces of controls[41].

252

253 In general, posterior estimates for the range of the spatial effects suggested small to medium

254 scale spatial correlation (Table 2; Figure 2). The range was smallest for groups A and H,

255 indicating that spatial correlation among points died off at relatively short distances. Both the

256 range and variance were largest for group C, indicating the spatial effect spanned a greater

257    distance but also displayed 'peaks'. This could be due to the presence of two distinct areas of

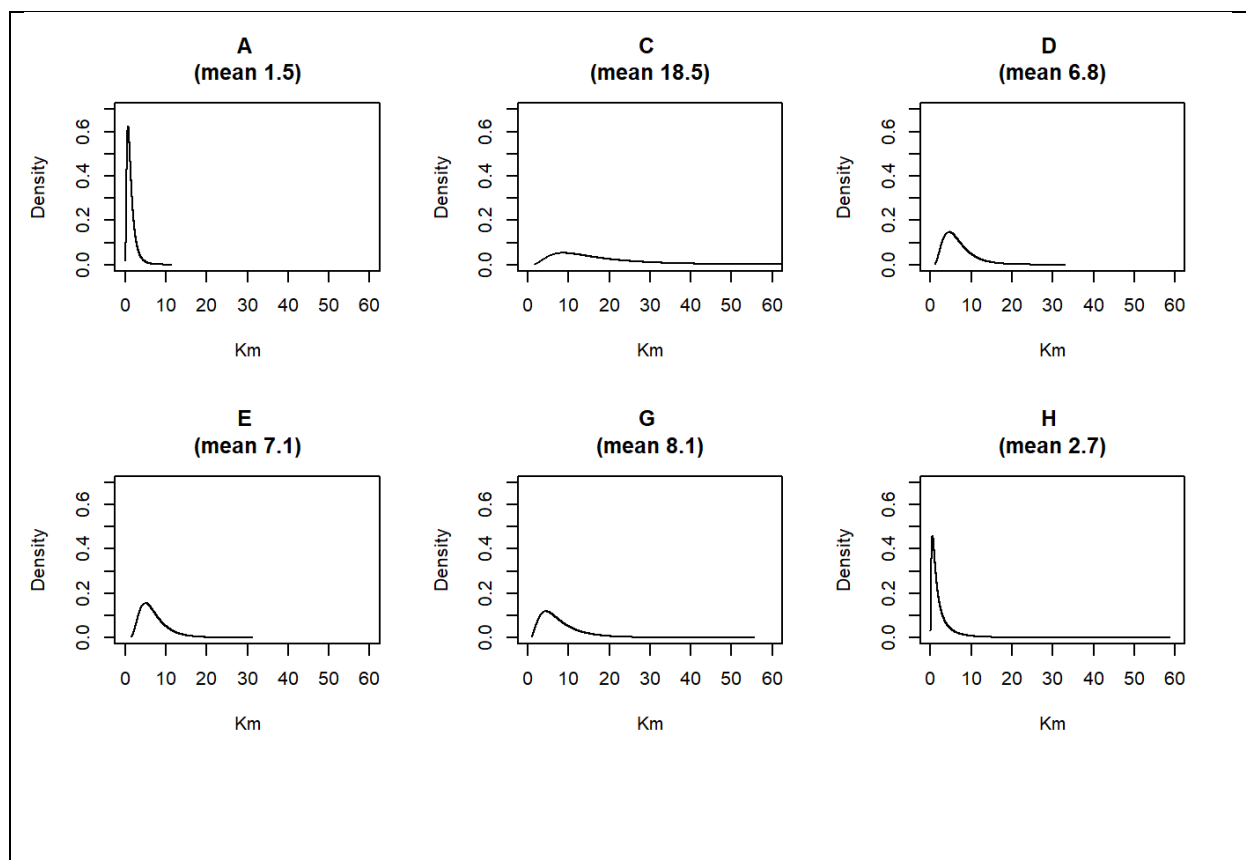258    increased intensity located relatively far from one another.

259

260    Table 2. Posterior mean estimates of the range and variance of the spatial effect for each outbreak

261    group (A-H), Gaborone, Botswana, 2012-2016.

|  | A | C | D | E | G | H |
|---|---|---|---|---|---|---|
| Range (mean) | 1.5 | 18.5 | 6.8 | 7.1 | 8.1 | 2.7 |
| Variance (mean) | 0.1 | 2.0 | 0.6 | 1.4 | 1.8 | 0.5 |

262

263    Figure 2. Posterior mean estimates and marginal distributions of the range and variance of the

264    estimated spatial effect for each outbreak group, Gaborone, Botswana, 2012-2016

Posterior mean and marginal distribution for range for spatial effect (outbreak groups A – H)
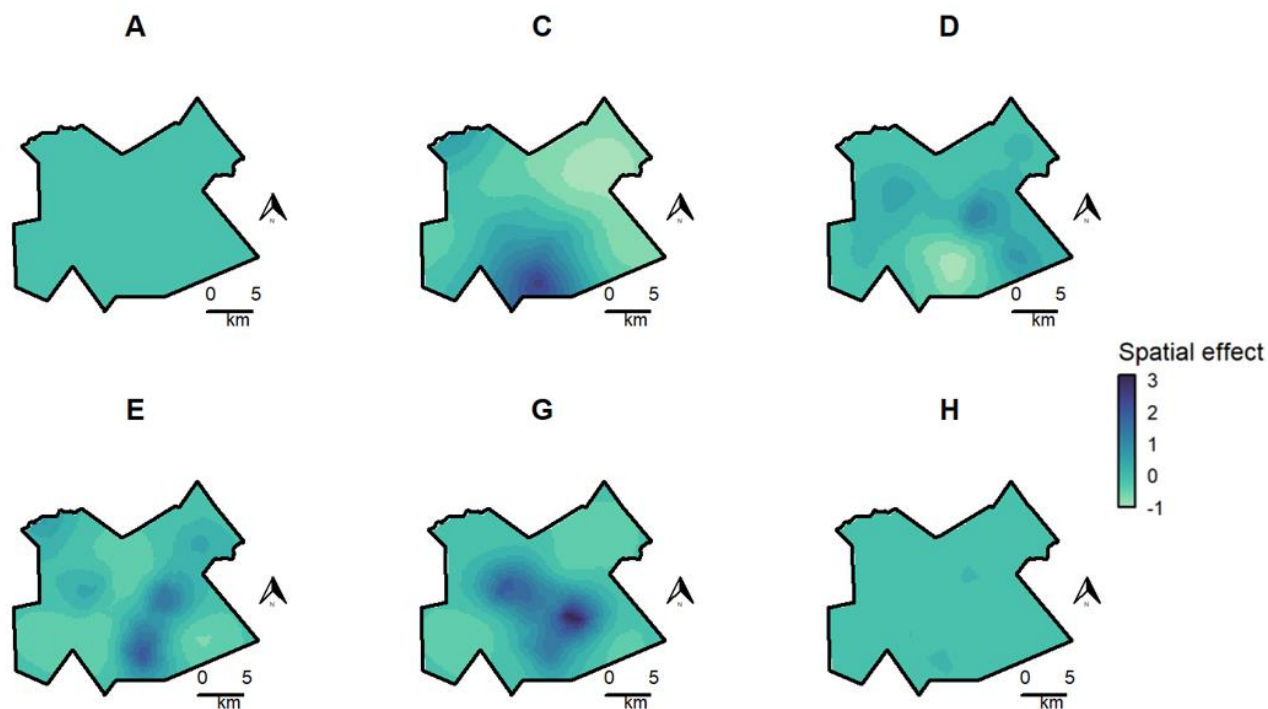
265

Maps of posterior mean estimates of spatial effects displayed different spatial patterns for each

outbreak group (Figure 3). Estimated spatial effects for group A and group H showed several

relatively small and dispersed areas of increased values compared to controls. Group C had a

notable area of increased spatial effect in the central southern part of the study area. Group D and

group G had two to three main areas of increased values that followed a broad east-west spread,

while for group E areas of increased estimates had a general north-south configuration.


Results of the sensitivity analysis using subsets of 70 and 140 randomly selected controls found

very similar results in terms of the spatial patterns and magnitude of estimated spatial effect by

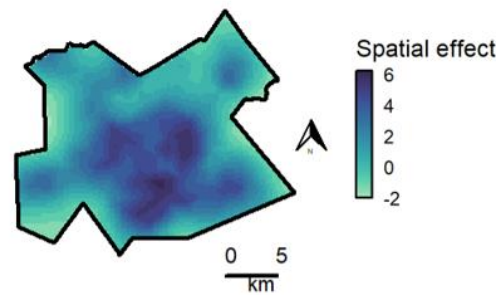group (Supplementary Figure 2 and Supplementary Figure 4).

277    Figure 3. Posterior mean estimates of spatial random effect for each outbreak group (A-H) and

278    controls (ungrouped participants), Gaborone, Botswana, 2012-2016. Values are shown on the

279    internal linear predictor scale and represent the contribution of the spatial random effect on the

280    response, after accounting for other fixed and random effects in the model. Departures from

281    baseline (above or below zero) for outbreak groups measure group-specific spatial patterns that

282    are not accounted for by the spatial distribution of activity spaces of controls. Darker colors

283    correspond to increased spatial effect estimates. Values are displayed on the same color scale for

284    all outbreak groups, though on a separate color scale for controls due to difference in sample

285    size.

286



287

288

289

290     Predicted spatial intensity

291     Maps of predicted mean spatial intensity displayed unique spatial patterns for each outbreak

292     group (Figure 4).

293

294     For group A, areas of increased spatial intensity of activity spaces followed an overall similar

295     pattern as seen for controls, with areas of highest intensity toward the center of the study area.

296     Group C had a distinct area of high intensity in the central southern part of the study area. Group

297     D had a notable area of increased intensity in the central east part of the map. Areas of highest

298     intensity for group E were in the central and south east, and for group G in the central east. The

299     areas of highest intensity for group H were located toward the center of the study area and also

300     resembled the overall spatial pattern seen for controls, though predicted values were relatively
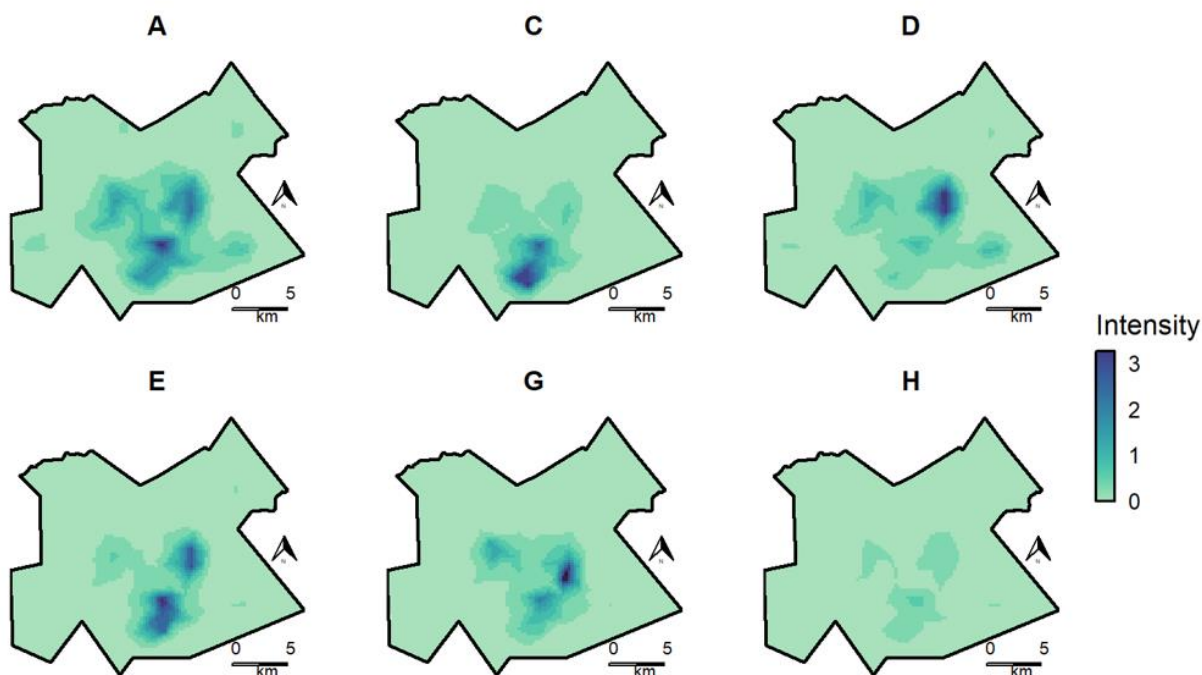
301     small compared to the other groups and not easily visible when mapped on the same color scale.
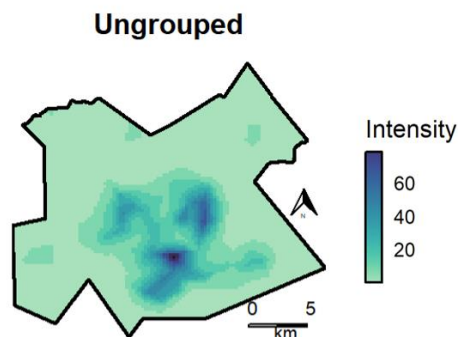
302

303     Results of the sensitivity analysis using subsets of 70 and 140 randomly selected controls found

304     very similar results for predicted spatial intensity by outbreak group (Supplementary Figure 3

305     and Supplementary Figure 5).

306

307     Figure 4. Predicted mean spatial intensity of activity spaces for participants in each outbreak

308     group (A-H) and controls (ungrouped participants), Gaborone, Botswana, 2012-2016. Values are

309     displayed on the response scale (obtained by exponentiating the linear predictor) and represent

310     predicted numbers of activity spaces per unit area (approximately 0.25 x 0.25 km). Areas of

311     increased intensity correspond to higher geographic concentration of activity spaces for

312     participants in each group. Intensity values are displayed on the same color scale for all outbreak

313     groups for ease of visual comparison, though on a separate color scale for controls due to

314     difference in sample size.



315

316

317

318    Exceedance maps

319    Exceedance maps for estimated spatial effects showed areas where the posterior mean had a high

320    probability (0.95) of being above 0 (greater than baseline) (Figure 5). Areas of significantly

321    increased spatial effect estimates based on full activity space analysis were detected for groups

322    C, D, E, and G. Groups A and H did not have areas meeting the specified threshold, which may

323    be due to a spatial distribution of activity spaces that resembles that of the control group.

324

325    Areas of significantly increased spatial effect estimates based on residential location alone were

326    detected for groups Cand E, though not for groups A, D, G, or H. For group E, exceedance areas

327    based on activity space were geographically broader than those based solely on residential

328    location. For group C, exceedance areas were similar in both analyses.

329

330    Figure 5: Exceedance maps for spatial effect greater than 0 (departure above baseline) with high

331    probability (0.95) – full activity space and residential locations only, Gaborone, Botswana, 2012-

332    2016.

5.1 Full activity space

## 5.2 Residential only



333

Exceedance maps for predicted spatial intensity values showed distinct areas of high spatial

intensity ('core areas') for each outbreak group, corresponding to areas where posterior mean

intensity values had a high probability (0.95) of being in the upper ten percent of estimates for

that group (Figure 6).

338

339     In general, core areas based on residential location alone were geographically restricted

340     compared to core areas based on full activity space analysis. For all groups core areas based on

341     full activity space analysis were larger than those based on residential locations alone. For groups

342     A, D, G, and H, core areas based on activity space also involved additional geographic focal

343     areas.

344

345     Figure 6. Exceedance maps for predicted spatial intensity to display 'core areas' by outbreak

346     group based on full activity space and residential locations only, Gaborone, Botswana, 2012-

347     2016

### 6.1 Full activity space

## 6.2 Residential only



348

**Discussion**

350    In our analysis, we detected geographically distinct patterns of activity space associated with

351    different TB outbreak groups. Core areas ('hotspots') of highest spatial concentration of activity

352    spaces for each group were located in different areas, with some being more geographically

353    widespread and others more compact. For outbreak groups C, D, E, and G, we detected areas

354    where the spatial concentration of activity spaces of grouped participants was significantly

355    higher than the baseline spatial distribution of activity spaces of ungrouped controls (increased

356    spatial effect). This could suggest that distinct areas of localized transmission play an important

357    role in these outbreaks. The spatial distribution of activity spaces for groups A (the largest

358    outbreak group) and H (the smallest) resembled the overall spatial distribution of activity spaces

359    belonging to the control group. The differences in spatial characteristics among the groups could

360    potentially correspond to the timing of how long a genotype of TB has been circulating in the

361    community. It could also represent transmission among socially or geographically distinct

362    contact networks.

363

364    Activity space hotspots could represent potential high-priority areas for spatially targeted

365    interventions such as active case finding for TB and other infectious diseases. This may be

366    particularly useful for outbreaks involving localized transmission. Exceedance maps displaying

367    core areas of spatial intensity and areas of increased spatial effects such as those shown above

368    could potentially be a useful tool for public health planning. We displayed static snapshots from

369    interactive maps at a relatively low spatial resolution to protect privacy, though in practice such

370    maps could be used to examine potential hotspots at different spatial scales.

371

372    We also found differences between exceedance areas detected using full activity space analysis

373    compared to residential location alone. Areas of core spatial intensity and significant spatial

374    effects based on activity space were generally larger, and sometimes included additional

375    geographic focal points, suggesting a notable portion of activity spaces may be located in areas

376    outside participants' home neighborhoods. A possible exception is group C, which had similar

377    exceedance areas in both analyses, suggesting activity spaces for these participants may

378    generally be located in closer proximity to place of residence. Relatively high unemployment in

379    group C and fewer work locations may have contributed to this observation. Our results show

380   that analyses based on residential location alone may not fully represent the spatial

381   characterization of hotspots.

382

383   Spatial analysis for infectious disease transmission involves an inherent assumption that the

384   locations analyzed are important with regard to transmission. Activity space analysis

385   incorporates important locations in the community where TB transmission may occur, and may

386   reduce exposure misclassification and improve the geographic characterization of transmission

387   chains[42]. This has implications for planning and evaluating targeted interventions[43]. For example,

388   a recent TB modeling study found limited effectiveness of spatially targeted screening based on

389   proximity to household locations of people with incident TB in Peru[44]. However, hotspots of TB

390   incidence based on household locations do not necessarily correspond with hotspots of

391   transmission[43,44]. Activity space analysis may help address issues such as this in spatial analysis

392   for TB transmission.

393

394   Activity space analysis has an established history of use in fields such as social geography and

395   urban planning[45–47]. It fits naturally into the spatial epidemiology framework, which emphasizes

396   place and location-based health exposures[23,48]. This approach acknowledges space as a social

397   determinant of health and helps incorporate the influence of social and physical environments on

398   health outcomes[23,49]. However, there are relatively few examples of spatial analysis of activity

399   space in the TB literature. An early example that helped lay the foundation for activity space

400   analysis in TB research was a study to detect TB hotspots in Japan[50]. The study incorporated

401   spatial and genomic data, though at a relatively low resolution (spatial data were aggregated at

402 the census tract level and genotype clustering identified using IS6110-based restriction fragment

403 length polymorphism (IS6110-RFLP) analysis)[50].

404

405 Our results are in line with studies of TB in the US[25] and South Africa[21] that both noted

406 differences between 'high-risk' areas identified with density maps of activity spaces compared to

407 residential locations alone. These studies highlighted the potential importance of activity space

408 analysis, though neither study included genomic data.

409

410 Our results are also in line with a recent study in Peru that combined WGS and spatial data to

411 identify differences in activity spaces of genotypically related and unrelated cases and non-TB

412 controls[26]. Notably, this study highlighted the potential to draw on methodology used in spatial

413 ecology, such as using UDs to model activity space, both at the individual and group level[26]. The

414 approach taken in this study was to focus mainly on quantifying size (geographic area) and

415 amount of overlap among participants' UDs, rather than detecting specific high-risk areas in the

416 community. Our study expanded on these methods by using a spatial point process model, which

417 allowed us to incorporate measures of uncertainty and detect potential hotspots by identifying

418 high-confidence areas of highest spatial intensity.

419

420 A limitation of our study is that incorporating activity space may have resulted in including

421 locations that are not relevant to transmission. Another limitation is that we did not include

422 specific measures of temporality, which is also an important element of transmission dynamics.

423

424    Another limitation of this study is that we did not examine potential contributing risk factors

425    driving the observed spatial patterns. Spatial variation is often a proxy for the influence of

426    unmeasured variables that may include sociodemographic, social, structural, or environmental

427    factors impacting risk[36,51]. The spatial LGCP modeling approach can incorporate spatially-

428    referenced covariates[51] (such as population-level sociodemographic characteristics or

429    environmental variables); however these data were not available for our analysis. We focused

430    primarily on identifying geographic areas of increased risk, which could potentially be targeted

431    for outreach such as active case finding. However, further analysis could incorporate additional

432    data or modeling techniques to assess potential risk factors.

433

434    Another limitation of this study is an unknown number of missing cases, activity space locations,

435    and WGS data that could potentially alter geographic characterization of genotypic groups.

436    Although the original study had relatively high enrollment (4,331/5,515 persons diagnosed

437    during the study period), not every person with TB was included, such as those diagnosed but not

438    enrolled and cases that were not detected. In addition, the use of location data obtained through

439    patient interviews is subject to recall bias and underreporting[52]. Other methods of obtaining

440    location data, such as prospective GPS tracking, have been suggested as potential alternatives[26].

441    However, locations visited during the infectious period prior to diagnosis and study enrollment

442    were of primary interest in this context[53]. Further, a study comparing locations reported by

443    participants and locations captured via GPS loggers found that for three quarters of respondents,

444    over 70% of self-reported locations matched with the GPS data[54].

445

446    **Conclusion**

447    Integrated geospatial and genomic analysis of activity space may help identify potential high-risk

448    locations of sustained transmission in the community. Activity space analysis may improve the

449    geographic characterization of transmission 'hotspots' compared to analysis of residential

450    location alone. This could help with planning and mobilizing interventions to interrupt ongoing

451    transmission, and could provide a valuable tool for public health officials working to eliminate

452    TB among marginalized communities[7,10].

453

References

454  References

455  1. World Health Organization. Global Tuberculosis Report 2023. World Health Organization;

456  2023. Accessed March 3, 2024. https://www.who.int/publications-detail-redirect/9789240083851

457  2. Vesga JF, Hallett TB, Reid MJA, et al. Assessing tuberculosis control priorities in high-burden

458  settings: a modelling approach. Lancet Glob Health. 2019;7(5):e585-e595. doi:10.1016/S2214-

459  109X(19)30037-3

460  3. Auld SC, Shah NS, Cohen T, Martinson NA, Gandhi NR. Where is tuberculosis transmission

461  happening? Insights from the literature, new tools to study transmission and implications for the

462  elimination of tuberculosis. Respirol Carlton Vic. Published online June 5, 2018.

463  doi:10.1111/resp.13333

464  4. Shaweno D, Trauer JM, Doan TN, Denholm JT, McBryde ES. Geospatial clustering and

465  modelling provide policy guidance to distribute funding for active TB case finding in Ethiopia.

466  Epidemics. 2021;36:100470. doi:10.1016/j.epidem.2021.100470

467  5. Zetola NM, Moonan PK, Click E, et al. Population-based geospatial and molecular

468  epidemiologic study of tuberculosis transmission dynamics, Botswana, 2012–2016. Emerg Infect

469  Dis. 2021;27(3):835-844. doi:10.3201/eid2703.203840

470  6. Guthrie JL, Gardy JL. A brief primer on genomic epidemiology: lessons learned from

471  Mycobacterium tuberculosis. Ann N Y Acad Sci. 2017;1388(1):59-77.

472  doi:https://doi.org/10.1111/nyas.13273

473  7. Gardy JL, Loman NJ. Towards a genomics-informed, real-time, global pathogen surveillance

474  system. Nat Rev Genet. 2018;19(1):9-20. doi:10.1038/nrg.2017.88

475    8. Moonan PK, Ghosh S, Oeltmann JE, Kammerer JS, Cowan LS, Navin TR. Using Genotyping

476    and Geospatial Scanning to Estimate Recent Mycobacterium tuberculosis Transmission, United

477    States. Emerg Infect Dis. 2012;18(3):458-465. doi:10.3201/eid1803.111107

478    9. Shaweno D, Karmakar M, Alene KA, et al. Methods used in the spatial analysis of

479    tuberculosis epidemiology: a systematic review. BMC Med. 2018;16(1):193.

480    doi:10.1186/s12916-018-1178-4

481    10. Inzaule SC, Tessema SK, Kebede Y, Ogwell Ouma AE, Nkengasong JN. Genomic-informed

482    pathogen surveillance in Africa: opportunities and challenges. Lancet Infect Dis.

483    2021;21(9):e281-e289. doi:10.1016/S1473-3099(20)30939-7

484    11. Smith JP, Oeltmann JE, Hill AN, et al. Characterizing tuberculosis transmission dynamics in

485    high-burden urban and rural settings. Sci Rep. 2022;12(1):6780. doi:10.1038/s41598-022-10488-

486    2

487    12. Ribeiro FKC, Pan W, Bertolde A, et al. Genotypic and Spatial Analysis of Mycobacterium

488    tuberculosis Transmission in a High-Incidence Urban Setting. Clin Infect Dis Off Publ Infect Dis

489    Soc Am. 2015;61(5):758-766. doi:10.1093/cid/civ365

490    13. Zelner JL, Murray MB, Becerra MC, et al. Identifying Hotspots of Multidrug-Resistant

491    Tuberculosis Transmission Using Spatial and Molecular Genetic Data. J Infect Dis.

492    2016;213(2):287-294. doi:10.1093/infdis/jiv387

493    14. Li M, Lu L, Jiang Q, et al. Genotypic and spatial analysis of transmission dynamics of

494    tuberculosis in Shanghai, China: a 10-year prospective population-based surveillance study.

495    Lancet Reg Health West Pac. 2023;38:100833. doi:10.1016/j.lanwpc.2023.100833

496    15. Moonan PK, Oppong J, Sahbazian B, et al. What Is the Outcome of Targeted Tuberculosis

497    Screening Based on Universal Genotyping and Location? Am J Respir Crit Care Med.

498    2006;174(5):599-604. doi:10.1164/rccm.200512-1977OC

499    16. Shaweno D, Trauer JM, Doan TN, Denholm JT, McBryde ES. (Pre)Geospatial clustering and

500    modelling provide policy guidance to distribute funding for active TB case finding in Ethiopia.

501    Epidemics. Published online May 19, 2021:100470. doi:10.1016/j.epidem.2021.100470

502    17. Dowdy DW, Golub JE, Chaisson RE, Saraceni V. Heterogeneity in tuberculosis transmission

503    and the role of geographic hotspots in propagating epidemics. Proc Natl Acad Sci U S A.

504    2012;109(24):9557-9562. doi:10.1073/pnas.1203517109

505    18. Shrestha S, Reja M, Gomes I, et al. Quantifying geographic heterogeneity in TB incidence

506    and the potential impact of geographically targeted interventions in south and north city

507    corporations of Dhaka, Bangladesh: a model-based study. Epidemiol Infect. Published online

508    April 19, 2021:1-27. doi:10.1017/S0950268821000832

509    19. Reid MJA, Arinaminpathy N, Bloom A, et al. Building a tuberculosis-free world: The Lancet

510    Commission on tuberculosis. The Lancet. 2019;393(10178):1331-1384. doi:10.1016/S0140-

511    6736(19)30024-8

512    20. Nelson KN, Shah NS, Mathema B, et al. Spatial Patterns of Extensively Drug-Resistant

513    Tuberculosis Transmission in KwaZulu-Natal, South Africa. J Infect Dis. 2018;218(12):1964-

514    1973. doi:10.1093/infdis/jiy394

515    21. Peterson ML, Gandhi NR, Clennon J, et al. Extensively drug-resistant tuberculosis hotspots

516    and sociodemographic associations in Durban, South Africa. Int J Tuberc Lung Dis Off J Int

517    Union Tuberc Lung Dis. 2019;23(6):720-727. doi:10.5588/ijtld.18.0575

518    22. Yang C, Lu L, Warren JL, et al. Internal migration and transmission dynamics of tuberculosis

519    in Shanghai, China: an epidemiological, spatial, genomic analysis. Lancet Infect Dis.

520    2018;18(7):788-795. doi:10.1016/S1473-3099(18)30218-4

521    23. Kestens Y, Wasfi R, Naud A, Chaix B. "Contextualizing Context": Reconciling

522    Environmental Exposures, Social Networks, and Location Preferences in Health Research. Curr

523    Environ Health Rep. 2017;4(1):51-60. doi:10.1007/s40572-017-0121-8

524    24. Matthews SA, Yang TC. Spatial Polygamy and Contextual Exposures (SPACEs): Promoting

525    Activity Space Approaches in Research on Place and Health. Am Behav Sci. 2013;57(8):1057-

526    1081. doi:10.1177/0002764213487345

527    25. Worrell MC, Kramer M, Yamin A, Ray SM, Goswami ND. Use of Activity Space in a

528    Tuberculosis Outbreak: Bringing Homeless Persons Into Spatial Analyses. Open Forum Infect

529    Dis. 2017;4(1):ofw280. doi:10.1093/ofid/ofw280

530    26. Bui DP, Chandran SS, Oren E, et al. Community transmission of multidrug-resistant

531    tuberculosis is associated with activity space overlap in Lima, Peru. BMC Infect Dis.

532    2021;21(1):275. doi:10.1186/s12879-021-05953-8

533    27. Baker CR, Barilar I, de Araujo LS, et al. Use of High-Resolution Geospatial and Genomic

534    Data to Characterize Recent Tuberculosis Transmission, Botswana. Emerg Infect Dis.

535    2023;29(5):977-987. doi:10.3201/eid2905.220796

536    28. Zetola NM, Modongo C, Moonan PK, et al. Protocol for a population-based molecular

537    epidemiology study of tuberculosis transmission in a high HIV-burden setting: the Botswana

538    Kopanyo study. BMJ Open. 2016;6(5). doi:10.1136/bmjopen-2015-010046

539    29. ESRI. ArcGIS Desktop. Published online 2019.

540    30. Diggle PJ. Statistical Analysis of Spatial and Spatio-Temporal Point Patterns. 0 ed. Chapman

541    and Hall/CRC; 2013. doi:10.1201/b15326

542    31. Banerjee S, Carlin BP, Gelfand AE. Hierarchical Modeling and Analysis for Spatial Data.

543    CRC Press; 2014.

544    32. Watson J, Joy R, Tollit D, Thornton SJ, Auger-Méthé M. Estimating animal utilization

545    distributions from multiple data types: A joint spatiotemporal point process framework. Ann

546    Appl Stat. 2021;15(4). doi:10.1214/21-AOAS1472

547    33. Simpson D, Illian JB, Lindgren F, Sørbye SH, Rue H. Going off grid : computationally

548    efficient inference for log-Gaussian Cox processes. Published online March 2016.

549    doi:10.1093/biomet/asv064

550    34. Lindgren F, Rue H. Bayesian Spatial Modelling with R-INLA. J Stat Softw. 2015;63(19):1-

551    25.

552    35. Illian JB, Sørbye SH, Rue H. A toolbox for fitting complex spatial point process models

553    using integrated nested Laplace approximation (INLA). Ann Appl Stat. 2012;6(4):1499-1530.

554    doi:10.1214/11-AOAS530

555    36. Moraga P. Spatial Statistics for Data Science: Theory and Practice with R. Chapman &

556    Hall/CRC Data Science Series; 2023. Accessed November 5, 2023.

557    https://www.paulamoraga.com/book-spatial/index.html

558    37. Krainski, Gómez-Rubio, Bakka, et al. Advanced Spatial Modeling with Stochastic Partial

559    Differential Equations Using R and INLA.; 2019. Accessed November 5, 2023.

560    https://becarioprecario.bitbucket.io/spde-gitbook/index.html

561    38. Gómez-Rubio V. Bayesian Inference with INLA.; 2021. Accessed November 7, 2023.

562    http://becarioprecario.bitbucket.io/inla-gitbook/index.html

563    39. Rue H, Martino S, Chopin N. Approximate Bayesian inference for latent Gaussian models by

564    using integrated nested Laplace approximations. J R Stat Soc Ser B Stat Methodol.

565    2009;71(2):319-392. doi:10.1111/j.1467-9868.2008.00700.x

566    40. Bachl FE, Lindgren F, Borchers DL, Illian JB. inlabru: an R package for Bayesian spatial

567    modelling from ecological survey data. Methods Ecol Evol. 2019;10(6):760-766.

568    doi:10.1111/2041-210X.13168

569    41. Palmí-Perales F, Gómez-Rubio V, López-Abente G, Ramis R, Sanz-Anquela JM, Fernández-

570    Navarro P. Approximate Bayesian inference for multivariate point pattern analysis in disease

571    mapping. Biom J. 2021;63(3):632-649. doi:10.1002/bimj.201900396

572    42. Keshavjee S, Dowdy D, Swaminathan S. Stopping the body count: a comprehensive

573    approach to move towards zero tuberculosis deaths. The Lancet. 2015;386(10010):e46-e47.

574    doi:10.1016/S0140-6736(15)00320-7

575    43. Huang CC, Trevisi L, Becerra MC, et al. Spatial scale of tuberculosis transmission in Lima,

576    Peru. Proc Natl Acad Sci U S A. 2022;119(45):e2207022119. doi:10.1073/pnas.2207022119

577    44. Havumaki J, Warren JL, Zelner J, et al. Spatially-targeted tuberculosis screening has limited

578    impact beyond household contact tracing in Lima, Peru: A model-based analysis. PLOS ONE.

579    2023;18(10):e0293519. doi:10.1371/journal.pone.0293519

580    45. Browning CR, Soller B. Moving Beyond Neighborhood: Activity Spaces and Ecological

581    Networks As Contexts for Youth Development. Cityscape Wash DC. 2014;16(1):165-196.

582    46. Horton FE, Reynolds DR. Effects of Urban Spatial Structure on Individual Behavior. Econ

583    Geogr. 1971;47(1):36. doi:10.2307/143224

584    47. Xi W, Calder CA, Browning CR. Beyond Activity Space: Detecting Communities in

585    Ecological Networks. Ann Am Assoc Geogr. 2020;110(6):1787-1806.

586    doi:10.1080/24694452.2020.1715779

587    48. Elliott P, Wartenberg D. Spatial Epidemiology: Current Approaches and Future Challenges.

588    Environ Health Perspect. 2004;112(9):998-1006. doi:10.1289/ehp.6735

589    49. Ortblad KF, Salomon JA, Bärnighausen T, Atun R. Stopping tuberculosis: a biosocial model

590    for sustainable development. Lancet Lond Engl. 2015;386(10010):2354-2362.

591    doi:10.1016/S0140-6736(15)00324-4

592    50. Izumi K, Ohkado A, Uchimura K, et al. Detection of Tuberculosis Infection Hotspots Using

593    Activity Spaces Based Spatial Approach in an Urban Tokyo, from 2003 to 2011. PLoS ONE.

594    2015;10(9). doi:10.1371/journal.pone.0138831

595    51. Diggle PJ, Moraga P, Rowlingson B, Taylor BM. Spatial and Spatio-Temporal Log-Gaussian

596    Cox Processes: Extending the Geostatistical Paradigm. Stat Sci. 2013;28(4). doi:10.1214/13-

597    STS441

598    52. Surie D, Fane O, Finlay A, et al. Molecular, Spatial, and Field Epidemiology Suggesting TB

599    Transmission in Community, Not Hospital, Gaborone, Botswana. Emerg Infect Dis.

600    2017;23(3):487-490. doi:10.3201/eid2303.161183

601    53. Bui DP, Oren E, Roe DJ, et al. A Case-Control Study to Identify Community Venues

602    Associated with Genetically-clustered, Multidrug-resistant Tuberculosis Disease in Lima, Peru.

603    Clin Infect Dis. 2019;68(9):1547-1555. doi:10.1093/cid/ciy746

604    54. Kestens Y, Thierry B, Shareck M, Steinmetz-Wood M, Chaix B. Integrating activity spaces in

605    health research: Comparing the VERITAS activity space questionnaire with 7-day GPS tracking

606    and prompted recall. Spat Spatio-Temporal Epidemiol. 2018;25:1-9.

607    doi:10.1016/j.sste.2017.12.003

608