

1 **Microbiota composition is moderately associated with**
2 **greenspace composition in a UK cohort of Twins.**

3

4 Authors:

- 5 • Ruth C. E. Bowyer¹ (ruth.c.bowyer@kcl.ac.uk)
- 6 • Caoimhe Twohig-Bennett ² (caoimhe.bennett@eruk.org.uk)
- 7 • Emma Coombes² (emma.coombes@uea.ac.uk)
- 8 • Philippa M. Wells¹ (philippa.wells@kcl.ac.uk)
- 9 • Tim Spector¹ (tim.spector@kcl.ac.uk)
- 10 • Andy P. Jones² (A.P.Jones@uea.ac.uk)[‡]
- 11 • Claire J. Steves^{1,3} [±] * (claire.j.steves@kcl.ac.uk)

12

13 [‡]AP Jones and CJ Steves contributed equally

14 *Corresponding author

15

16 Affiliations:

17 1. Department of Twin Research and Genetic Epidemiology, King's College London, 3-4th

18 Floor South Wing Block D, St Thomas' Hospital, Westminster Bridge Road, London SE1

19 7EH, UK

20 2. Norwich Medical School, University of East Anglia, Norwich, Norfolk, NR4 7TJ

21 3. Department of Ageing and Health, St Thomas' Hospital, 9th floor, North Wing,

22 Westminster Bridge Road, London SE1 7EH

23

24 **Abstract**

25 Exposure to natural environments, known as greenspace, appears to positively influence
26 health, yet the mechanisms are unclear. Given that gut microbiota are associated with
27 inflammatory disorders more prevalent in urban areas and individuals with lower greenspace
28 exposure, microbiota may act as a mediator between greenspace and health. Using 2443
29 participants of the TwinsUK cohort, microbiota differences were compared in relation to
30 rural/urban living and with quantiles of area-level greenspace at three different
31 neighbourhood distances: 800m, 3000m and 5000m. Using microbiota data captured from
32 faecal samples using 16S rRNA marker gene sequencing, small compositional differences in
33 association with 3000m greenspace ($p=0.003$) in models adjusted for confounders of
34 microbiota variance (sequencing depth, antibiotics use, body mass index, frailty, age, diet,
35 region and socioeconomic variables) were observed. Differences in abundances of genus
36 were observed for all measures of greenspace in adjusted models; a key pathogenic genus
37 was increased in abundance in association with urbanicity (*Escherichia/Shigella*, $\log_{FC} =$
38 0.73742 , $p_{adj} < 0.001$). Further, utilising the twin structure, within-pair differences in
39 microbiota composition were compared and associations with 800m greenspace observed
40 (factor level significance in association with greatest difference, $\beta = 0.08$, $p = 0.0162$) as were
41 differences in *Escherichia/Shigella*. The microbiota signature of those with a greater
42 exposure to greenspace, but not necessarily explicitly rural individuals, was distinct from
43 other individuals, suggesting microbiota as a potential mediator for greenspace and health.

44

45 **Keywords: Greenspace, microbiota, rural-urban classification, twin differences**

46

47 **Introduction**

48

49 The human gut microbiome is the collective genome and surrounding environment of the assemblage
50 of microorganisms inhabiting the gastrointestinal tracts (Goodrich et al., 2016; Marchesi and Ravel,
51 2015). It interacts with the immune system, is essential for nutrient processing, and offers protection
52 from pathogens (Geva-Zatorsky et al., 2017; Litvak and Bäumlér, 2019; Magnúsdóttir et al., 2015).
53 The microbiome is therefore increasingly considered a target for clinical intervention. Whilst the
54 extent of its influence on health is still being debated, there is evidence that it can be considered a
55 biomarker and potential mediator of environmental exposures and health (Cresci and Bawden, 2015).
56 Indeed, because of the relatively low influence of genetics on its composition, environmental factors
57 are key to understanding the forces shaping its composition (Rothschild et al., 2018). Whilst
58 influences such as diet, medication use and disease status have been more widely studied (for
59 example: Ordiz *et al.*, 2015; Shreiner, Kao and Young, 2015; Jackson *et al.*, 2016; Belzer *et al.*, 2017;
60 Weersma, Zhernakova and Fu, 2020), the influence of environmental factors more distal to the host
61 on microbiota composition, such as where an individual lives, are less well understood (Schmidt et al.,
62 2018). This is despite both the emerging importance of the microbiome and increasing evidence that
63 exposures from the environment influence health.

64 High levels of exposure to natural environments may, for certain conditions, have similar therapeutic
65 benefits equivalent in magnitude to some drug treatments (Twohig-Bennett and Jones, 2018) and
66 consequently, urban greenspaces have been proposed as a means to mitigate the negative health
67 impacts of urbanisation. Whilst frequently proposed because of the association of inflammatory and
68 chronic immunological disease and urbanisation, it is unclear whether the mode of action for the
69 observed benefit is mediated by the gut microbiota (Rook, 2013; Rook et al., 2013). Previous studies
70 have suggested differences in the microbiota by urbanicity (Obregon-Tito et al., 2015), but have not
71 assessed the relative importance of exposure to different types of environment.

72 Associations between any health trait and greenspace are likely confounded by the complex
73 interdependency of social and behavioural factors that contribute to an individual's greenspace
74 exposure (Robinson and Jorgensen, 2020; Stamper et al., 2016). This study makes use of the
75 TwinsUK cohort, the UK's largest registry of twins, and aims to investigate how microbiota
76 assemblages differ by greenspace or urbanicity.

77 **Methods**

78

79 **Participants**

80

81 Participants were members of TwinsUK, the UK's largest registry of mono and dizygotic twins, run
82 out of the Department of Twin Research & Genetic Epidemiology, King's College London. Incepted
83 to study the heritability of osteoarthritis and osteoporosis, the registry has grown to be one of the most
84 clinically detailed twin cohorts worldwide (Verdi et al., 2019). Data from twins is collected
85 approximately every four years either upon visit to the Clinical Research Facility, or via questionnaire
86 collected on a more frequent basis over the telephone, sent in the post, or filled in and brought to
87 visits. For the purposes of this analysis, all measures were matched based on date unless otherwise
88 stated. 3218 samples from 2707 individuals from within the TwinsUK cohort with appropriately
89 collected and processed 16S rRNA gene sequence data were considered for this analysis. 3010
90 samples were successfully geocoded; those where a postcode match could not be found were likely
91 due to an administrative error, recording of a now disused postcode, or because the participant lived
92 abroad. Participants were spread across the UK, with the majority in England (96.5%), of which
93 43.8% were living in South East of England or Greater London at time of microbiota sample. For
94 primary analysis, the earliest sample provided by individuals was used where applicable with repeat
95 samples provided by the same individual on the same day removed (sample chosen at random for
96 removal), resulting in 2443 samples/individuals considered in this analysis.

97

98 **Microbiome samples**

99

100 Profiles of gut microbiota composition for each twin were available as a subset generated from faecal
101 samples as previously described (Goodrich et al., 2016). Samples were collected and sequenced on a
102 rolling basis between November 2011 and 2015. Participants stored samples in sealed ice packs and
103 either posted them to the research department or provided them during clinical visit. Samples were
104 stored at -80°C and shipped frozen to Cornell University for DNA extraction and amplification of
105 the V4 variable region of the 16S rRNA gene amplified in duplicate by PCR using 515F and 806R
106 primers, following isolation of genomic DNA using the 'MoBio PowerSoil htp DNA isolation kit'
107 from an aliquot of $\sim 100\text{mg}$ of each sample. As described by Goodrich and colleagues 2014, PCR
108 reactions comprised 2.5 U Easy-A high-fidelity enzyme, $1 \times$ buffer, 10-100 ng DNA template, and
109 $0.05 \mu\text{M}$ of each primer. Initial denaturation was carried out at 94°C for 3 min followed by 25 cycles
110 of denaturation at 94°C for 45 seconds, followed by annealing at 50°C for 60 seconds, extension at
111 72°C for 90 seconds, and a final extension at 72°C for 10 minutes. A magnetic bead system was used
112 to combine and purify the PCR reactions and washout contaminants. The QuantiT PicoGreen dsDNA
113 Assay Kit was used to quantify the PCR amplicons, aliquots of which were combined for a final
114 concentration of approximately $15 \text{ ng}/\mu$. This 16S rRNA gene is a unit of the bacterial genome that
115 exhibits extreme sequence conservation whilst containing 9 hypervariable regions that can be targeted
116 for broad differentiation between microbial lineages (Tringe and Hugenholtz, 2008). For the
117 limitations and advantages to this approach see Weinstock (2012) and Pollock et al. (2018).

118

119 The resulting sequences were analysed as amplicon sequence variants (ASVs) following the DADA2
120 pipeline (Callahan et al., 2016). Briefly, DNA sequences were demultiplexed, and separate forward
121 and reverse read files were generated for each sample using QIIME (Caporaso et al., 2010). Quality of
122 sequences was assessed, with ends trimmed to remove poor quality reads, error estimated within-
123 sample for forward and reverse reads, and then the ASV algorithm applied. Forward and reverse
124 ASVs were joined, and the total dataset merged. Chimeras were removed. Taxonomic assignment was
125 via SILVA 1.3.2. Samples with less than 10,000 sequences or with only one viable read direction

126 were removed. ASVs can be considered as a ‘taxonomic unit’ of the microbiome, akin to a species of
127 microorganism as best estimated using the 16S rRNA gene sequencing method, with taxonomic
128 description at the lowest resolution possible for that sequence (i.e. some ASVs will be mapped to a
129 species, others will only be assigned to a taxonomic order or family).

130 **Residential location**

131

132 To maximise the likelihood that the participant was living at the address when the sample was
133 provided, residential postcodes (zipcodes) of study participants were linked to microbiota samples by
134 the nearest date of a provided postcode prior to the sample date rather than the date outright. Matched
135 postcodes were then geocoded based on their centroids (centre points) in R using package
136 ‘PostcodesioR’ (Walczak, 2019), a wrapper that allows R interface with Postcodes.io and maps
137 postcodes to Ordnance Survey postcode centroids.

138

139 **Key predictor variables**

140

141 The environmental exposure variables considered in this study were measures of urbanicity and
142 greenspace, computed as follows.

143

144 To measure urbanicity, the Rural-Urban 2-fold classification (RUC) was linked to each individual’s
145 postcode as a measure of the urbanicity of their local area. The RUC was defined at the 2011 UK
146 Census ‘output-area’ level (For data sources, please refer to *data availability*). The measure was
147 included to compare whether simply urbanicity, rather than greenspace, might shape microbiota
148 composition.

149

150 For the measurement of greenspace in the residential vicinity of each study participant, buffer zones
151 around each postcode were calculated using the Ordnance Survey Open Roads dataset that indicates
152 the footpath and roads network across the UK. Using the QGIS 3.4.2 package, the area covering the

153 800m, 3000m and 5000m road and path networks from postcode centroids was converted to a
154 polygon (closed geographical unit) representing the area surrounding an individual's home based on
155 travel distance. Buffer distances were chosen to represent the effect of the immediate environment
156 surrounding the home as well as the wider environment an individual might move through on a daily
157 basis (Hillsdon et al., 2015).

158

159 The 25m² Land Cover Map of Great Britain 2015 (LCM, version 1.2) was downloaded from the
160 Centre for Ecology & Hydrology via the 'Digimap' data portal. The LCM uses satellite data to
161 classify landcover across the UK into 21 classes. Buffer polygons were overlaid with the LCM, and
162 percentage of each landcover type falling within each was calculated, with all the 21 classes that were
163 not considered as 'urban' or 'suburban' being considered as greenspace. For the purposes of this
164 analysis, greenspace percentages were defined as within-population quantiles. For details on the data
165 sources please refer to the data availability section.

166

167 **Statistical analysis**

168

169 All statistical analyses were undertaken in RStudio using R version 3.6.3. Missing covariate data
170 (Body Mass Index BMI, frailty, diet and antibiotic use) were imputed using the 'missForest' R
171 package, an iterative imputation method based on random forest, that can impute both continuous and
172 categorical data (Stekhoven and Bühlmann, 2012). Covariates were compared for their differences by
173 the key outcome variables (urbanicity and greenspace) using pairwise Wilcoxon rank sum tests or chi-
174 squared tests depending on the variable. Means, standard deviations (for continuous variables) and
175 percentage of the dominant class (factor variables) were calculated.

176 *Compositional (beta diversity) analysis*

177 A compositional framework was favoured for comparison of inter-individual community composition
178 (sometimes referred to as 'Beta diversity') (Gloor et al., 2017). Hypothesising that the composition

179 and presence of rarer microbes were more likely to differ geographically than ubiquitous ones, and
180 equally, that it was of interest to understand whether the composition of the conserved microbial
181 function differed spatially, ASVs were subset to those ‘abundant’ (>20% of the population, a
182 frequently used cut off in microbiota analysis) and ‘rare’ (ASVs in between 7.5-20% of the
183 population). As a noise-reduction strategy, abundant ASV table were ‘collapsed’ to genus (i.e.
184 transformed counts were grouped together where their taxonomic assignment suggested they belonged
185 to the same genera).

186

187 For this aspect of the analysis, following methods of Tipton et al. (2018), a pseudocount was added to
188 the ASV table (prior to subsetting) using a Bayesian-multiplicative replacement of count zeros via the
189 ‘zCompositions’ package (Palarea-Albaladejo and Martín-Fernández, 2015). Samples were subset as
190 described above (i.e. rare, abundant and genus), and relative abundance of the samples was calculated
191 in ‘phyloseq’ (McMurdie et al., 2013), followed by a centred-log ratio (CLR) transformation. Finally,
192 Aitchison’s distance calculated on each table.

193 Our analysis was framed by the use of permutational analysis of variance (PERMANOVA) using the
194 vegan package (Oksanen et al., 2015) to assess variation in microbiota composition in association
195 with outcome variables (i.e. how different are microbial assemblages in individuals living in areas of
196 high greenspace compared to low greenspace). The model works by comparing group with the null
197 hypothesis that group centroids and dispersion (defined by a distance matrix) are equivalent for all
198 groups (Anderson, 2001). Homogeneity of variance and pairwise comparison of factors were assessed
199 using *beta.dispers* and *pairwise.adonis* respectively.

200 PERMANOVA describes how variation of (in this case) a distance matrix can be attributed to
201 different factors and variables. When specifying the model, terms are added sequentially, and
202 therefore where factors are correlated, the extent of variance that they explain diminishes if they are
203 added after a corollary (Anderson, 2001). Therefore, a nested approach was taken, as follows -

204 **Model 1:** where models were performed with the habitat variable adjusted for family structure (to
205 account for twin relatedness) and library size (a technical variable representing the number of

206 sequences per sample). **Model 2:** models adjusted as in Model 1 with the addition of their geographic
207 region (to account for differences simply as a facet of distance/location, with the regions grouped as
208 follows: 1. East Midlands, 2. East of England, 3. London, 4. North England & Scotland, 5. South East,
209 6. South West, 7. West Midlands and 8. Yorkshire and the Humber). **Model 3:** A saturated model, as
210 in Model 2 with the inclusion of potential mediatory and confounding biological variables, all of
211 which have previously been shown to associate with microbiota composition in this cohort. They
212 were: body mass index (BMI), a measure of dietary quality in the form of the healthy eating index
213 (HEI) (Bowyer et al., 2018), the frailty index – a measure of health deficit derived using Rockwood’s
214 method (Searle et al., 2008), age at microbiota sample, antibiotic use in the previous month, highest
215 educational attainment and the area-based Index of Multiple Deprivation (IMD). The final two were
216 considered as key measures of socioeconomic status, with the latter being an important capture of
217 area-level deprivation. To understand the relative variance explained by each variable comparatively,
218 where statistically significant Model 3 was repeated to include all terms marginally.

219 As a result of the above methodology, three outcome measures were assessed in the sequential
220 manner as laid out above based on Aitchinson’s distance matrix calculated on abundant ASVs, rare
221 ASVs and genera, with the four key variables of interest as: 1) RUC, 2) 800m greenspace quantiles, 3)
222 3000m greenspace quantiles and 4) 5000m greenspace quantiles, totalling 36 models. Therefore the
223 multiple-testing alpha threshold was set as 0.004 for the PERMANOVAs, by assuming very
224 conservatively independence of the 3 Aitchinson’s distances and 4 outcome variables, but not
225 adjusting to account for the sequential design (i.e. $0.05/(3 \times 4)$).

226 *Differences in ASV and genera abundances*

227 Differences in ASV abundances were assessed using the DESeq2 pipeline (Ben J Callahan et al.,
228 2016; Love et al., 2014) using both the abundant ASVs and genera tables described above,
229 transformed using the variance stabilising method internal to DESeq2. Models were fitted both
230 unadjusted and adjusted for sequencing depth, antibiotics in the month prior to sample collection,
231 Body Mass Index, frailty, age, diet, index of multiple deprivation, highest educational attainment and

232 region of living. Results were false-discovery rate adjusted using the the Benjamini–Hochberg
233 procedure.

234 *Twin-pair differences*

235 Twin pair differences were assessed in two ways. First, the within-in pair dissimilarity was extracted
236 from the ‘abundant’ distance matrix (described above) and used as the outcome variable in linear
237 regressions, adjusting for within-pair difference in antibiotic use, BMI, education attainment, frailty,
238 diet, index of multiple deprivation and library size. Difference was calculated as the absolute value
239 between each individual and their co-twin; all covariates are as described above. Second, variance
240 stabilised genus abundances (derived within the DESeq2 pipeline) for those statistically significant in
241 the population wide analysis were tested for difference between twins discordant for the relevant
242 factor of interest using paired Wilcox rank sum tests. Discordance was defined as where one twin
243 lived in a rural area versus one twin living in an urban area for RUC, and where one twin lived in a
244 different quantile of greenspace.

245

246 **Results**

247 From the 3218 samples from 2707 individuals considered for analysis, 3010 were successfully
248 geocoded and, after removal of repeats, 2433 samples were used in this analysis. Descriptive statistics
249 (**Table 1**) suggest broadly the same sub-populations were observed across each quantile and by
250 rural/urban differences, although there were some key differences: age was significantly lower in Q1
251 of 3000m greenspace and 5000m quantiles compared to all four other quantiles ($p < 0.001$) and to a
252 lesser extent in 800m greenspace quantiles (Q1:Q3 $p = 0.0022$, Q1:Q4 $p = 0.008$) and in urban
253 dwellers compared with rural ($p = 0.002$) and IMD differed by quantiles of 3000m and 5000m
254 greenspace, primarily the most deprived compared to the others (Q1:Q2, Q3, Q4 $p < 0.001$, Q2:Q3
255 $p = 0.01$, Q3:Q4 $p = 0.007$, Q1:Q2, Q3, Q4 $p < 0.01$, Q2:Q3 $p = 0.0003$, Q3:Q4 $p = 0.00125$).

256 ‘Out-of-bag’ imputation error (OOB error) was: 0.000169 (normalised root mean squared error,
257 NRMSE – continuous variables) and 0.0131 (Proportion of falsely classified- PFC – categorical
258 variables) suggesting the imputed dataset was suitable for analysis.

259 Percentage greenspace in buffers surrounding homes was highly correlated between the different
260 distances (**Figure S1**), although there were many more individuals with <5% estimated greenspace
261 within the 800m buffer (**Figure S2**).

262 *Microbiota differences are observable in association with greenspace and to a lesser extent*
263 *urbanicity*

264

265 Whilst there were statistically significant associations in unadjusted models for all greenspace buffers
266 and rare and abundant microbiota compositions, (**Table S3**) saturated models did not generally show
267 differences between any greenspace buffers nor urban-rural categories apart from a nominally
268 significant association with 800m and rare microbiota composition ($p=0.02$). However, when
269 considering genus differences, there were statistically significant differences observed according to
270 urban-rural status and across quartiles of all neighbourhood sizes (**Table 2**). Only 3000m greenspace
271 buffers showed statistically significant differences across quartiles in microbiota compositions in
272 saturated models, with post-hoc pairwise comparisons suggesting the association was primarily driven
273 by the fourth quantile representing the highest amount of greenspace. Comparison of the relative
274 variance explained by each included variable (marginal model) suggested that whilst overall variance
275 explained was low, greenspace composition explained nearly as much as diet (HEI) and more than the
276 included measures of socioeconomic status and health deficit (**Figure 1i**). Sensitivity analysis where
277 one twin was removed at random in each pair suggested diminished associations, with few models
278 passing the calculated multiple testing threshold of 0.004 (Sensitivity analysis S1).

279

280

281 ---- **Figure 1 here** -----

282

283 ***Differences in abundance of ASVs and genus were observed by urbanicity and greenspace***

284 Results of DESeq2 tests of differences of ASV (n=297) and genus (n=87) abundances suggested
285 several associations – particularly in adjusted models (Figure 1.ii-iii). The clearest signal was the
286 suggested increase in abundance of *Escherichia/Shigella* with urbanicity, which was replicated across
287 models of adjusted and unadjusted ASVs and genus (Genus, adjusted: logFC = 0.73742, padj <
288 0.001). Genera annotated as *Streptococcus* and *Prevotella_9* also increased in association with
289 urbanicity, whereas *Ruminiclostridium* also decreased in unadjusted models, but these signals were
290 not replicated in adjusted models. Only nominal significance was observed in unadjusted ASV models
291 of greenspace, whereas ASVs were associated with adjusted models of 3000m and 5000m greenspace
292 buffer, including one annotated at species level *Roseburia inulinivorans* with 3000m (logFC=
293 0.19951, padj=0.03957) and *Lachnoclostridium* with 5000m; the latter was replicated when
294 unadjusted genus models, but not adjusted models (Figure 1.iii-iii. Table S4). Other genera were
295 observed in association with all three greenspace buffers in adjusted models; *Haemophilus* decreased
296 in association with 800m greenspace.; *Escherichia/Shigella* decreased in association with 3000m
297 greenspace and *Ruminococcaceae_UCG_014* increased; *Anaerostipes* increased in association with
298 5000m greenspace (Figure 1.iii, Table S4)

299

300 ***Modest twin pair differences in greenspace and urbanicity reflect observations of difference in***
301 ***microbiota composition***

302 There were 975 twin pairs within the subset used for this analysis. Within-pair difference in
303 microbiota composition was associated with within-pair difference in 800m greenspace at factor level
304 (Twin pair regression, factor level significance in association with greatest difference, $\beta = 0.08$, $p =$
305 0.0162, but comparison with null model not statistically significant at $p = 0.055$), but not RUC or other
306 greenspace buffers. This suggests a small, statistically significant difference in microbiota comparison
307 between twins where one lives in an area of low greenspace immediate to their home compared with
308 the co-twin who lives in an area of high greenspace immediate to their home. The only other variable

309 that was statistically significant in these twin pair models apart from sequencing depth was within-
310 pair difference in IMD.

311 Of the five genus that were significant in the population-wide analysis, *Escherichia/Shigella*
312 abundances were statistically significantly different between twins discordant for rural-urban living
313 (Wilcox rank, paired, $V = 125$, $p\text{-value} = 0.02$) (Figure 2).

314

315 ---Figure 2 here---

316

317 Discussion

318

319 Subtle differences of microbiota in association with the environmental composition around an
320 individual's home postcode were observed. The microbiota signature of those with a greater
321 neighbourhood exposure to greenspace was distinct from other individuals, and a difference in ASV
322 and genus abundance was observed particularly between rural and urban dwelling individuals.

323 Differences in microbiota composition and geographic factors have previously been cited. In a study
324 of traditional populations in India, Dehingia and colleagues observed microbiota variation differed
325 within geographic region (Dehingia et al., 2016). He *et al.* note regional differences explained the
326 greatest variance in their study of microbiota divergence between districts of China (He et al., 2018).

327 In this study's PERMANOVA analysis, findings of He and colleagues were not replicated, instead
328 observing factors other than environmental influences to explain a greater variance of microbiota
329 composition. This may in part be explained by differences in how covariates were measured, such as
330 diet which was included as a single composite measure whereas He and colleagues included 19
331 different dietary 'items' which by themselves would individually capture less variance than a
332 composite. Reflecting this, and as observed previously in this data (Bowyer et al., 2019), adonis-R^2
333 values are low which suggests that whilst statistically significant, variation explained by all factors is
334 low. Indeed, because of the high inter-individual variation in the microbiota, measured factors in

335 general explain low proportions of microbiota variance and so should be considered relatively to one
336 another (Falony et al., 2016; Rothschild et al., 2018).

337 Compositional differences in rarer species were associated with local greenspace availability. This
338 could indicate an increase of exposure to rarer microbes in association with more greenspace
339 (although compositional differences assessed in this manner do not reflect directional differences), but
340 this was not reflected in subsequent twin-pair analysis, and importantly did not pass the multiple
341 testing threshold in saturated models. This suggests that other confounders could explain a degree of
342 variance associated with greenspace and microbiota. Area-level socioeconomic deprivation, included
343 in saturated models via the Index of Multiple Deprivation (IMD) could in part explain this as IMD
344 and greenspace are correlated, and has previously been cited as a confounder to studies of greenspace
345 and health. Unpicking this should be the focus of future work. Of particular interest in the twin-pair
346 models was that difference in IMD, along with 800m greenspace difference, was significantly
347 associated with twin dissimilarity. Importantly, being a composite area-based measure, the IMD does
348 not necessarily reflect the socioeconomic status of an individual living in an area (this is known as the
349 ‘ecological fallacy’) but rather the area-level socio-economics at the population level (i.e.
350 environmental deprivation). Considered together therefore, this strengthens the suggestion that the
351 physical environment might exert pressures that shape the microbiota.

352 Little compositional difference by explicitly measured urbanicity (as opposed to high greenspace
353 exposure) was observed, despite the existence of previous evidence to suggest the microbiota differs
354 between rural and urban environments. For example, a comparison between Russian cities and
355 villages observed differences in the associated microbiota, with the rural microbiome being more
356 health associated (Tyakht et al., 2014). The difference between urban and rural living in the UK
357 however may be small compared to the environment in that study. Nevertheless, the increase of
358 *Escherichia/Shigella* in association with urbanicity both in population-wide and twin pair models does
359 suggest subtle differences in microbiota between urban and rural individuals. *Escherichia* and
360 *Shigella* are genetically closely related pathogens (and not differentiable by a 16s rRNA marker gene
361 approach) and whilst these species are commonly observed as members of the commensal microbiota,

362 they are known pathogens and generally not associated with health. The increased colonisation in
363 urban individuals could suggest that the resistance to pathogens conferred by a healthy microbiota is
364 diminished in these areas.

365

366 The strengths of this study lie in the investigation of an understudied influencer of microbiota
367 variance using samples taken from a large and well characterised population cohort integrated with
368 environmental measures computed within a geographical information system, and the use of twin
369 structure to do so. There are however several limitations. Sensitivity analysis re-running
370 PERMANOVAs revealed that whilst nominal statistical significance remains and R^2 values are
371 similar, the results no longer pass our multiple testing threshold of 0.004. This may be due to the
372 reduction in numbers, but we cannot exclude that there might be an effect of twin relatedness. Where
373 the outcome of interest is only moderately heritable as in the case of the microbiota (Rothschild et al.,
374 2018) the twin design has limited power and thus analysis is paired with a population-wide approach.
375 Observations were weak and certain geographic variables of interest that may correlated with
376 greenspace and therefore act as confounders were not included. For example, air pollution could
377 promote negative health consequences by influencing the microbiota; exposure to nitro-polycyclic
378 aromatic hydrocarbons through diesel fuel emissions have been suggested to promote nitroreduction
379 by the bacteria that results in a microbiota dependent increase in carcinogenesis (Claus et al., 2016).
380 Whilst the complex interdependency between socioeconomics, health and environmental exposure is
381 partially mitigated via the use of twin pair and adjustment this may not fully account for these
382 potential confounders, and inferences are limited by use of static data to capture environmental
383 exposure. Dynamic data reflecting the behaviour of an individual, and direct measurement of their use
384 and time spent in greenspace would address some of the uncertainties raised here. A further limitation
385 was the inclusive definition of “greenspace” as being any area of non-urban environment, which
386 therefore does not consider factors such as the accessibility of land or the types of habitat present. A
387 final limitation is the use of 16S rRNA data, which whilst appropriate for initial observational studies
388 such as this, gives less insight than whole genome sequencing data which would better be able to

389 confirm the differences in microbial-induced immunoregulation between individuals via functional
390 annotation, and additionally allow capture of non-bacterial microorganisms.

391 **Conclusion**

392 This study has observed difference in the microbiota of twin pairs with different environmental
393 exposures, suggesting elements of the host's habitat may contribute to the complex compositions of
394 microbial assemblages. Considered together, results are suggestive that there are geographic patterns
395 in the microbiota observable in this dataset which do not seem to be accounted for by diet, BMI and
396 frailty. Modest microbiota variance may be attributable to greenspace within a close to moderate
397 distance to an individual's home and urbanicity. Findings highlight the potential importance of
398 considering non-lifestyle factors that influence the composition of the microbiota and add to the
399 literature exploring the microbiome-environment-health axis.

400

401

402

403

404

405 **Figure legends**

406

407 **Figure 1.i.** Comparison of adonis-R² values from saturated Genera PERMANOVA (adonis) marginal
408 model, demonstrating the comparative variance of each variable in a saturated model. Stars indicate
409 statistical significance at 2000 permutations: ***= p<0.0001, ** = p < 0.001, *= p< 0.004. ii.

410 Volcano plots of differences in abundance of ASVs adjusted for sample sequence depth, antibiotics in
411 the month prior to sample collection, Body Mass Index, frailty, age, diet, index of multiple
412 deprivation, highest educational attainment and region. Log fold change (logFC) represent change in
413 relation to A. Urban living, B. Higher greenspace percentage in 800m, C. Higher greenspace
414 percentage in 3000m and D. Higher greenspace percentage in 5000m. Grey point = non-significant
415 (n.s), blue point = nominal significance (p<0.05) and red point = fdr-adjusted significance (q<0.05)
416 **iii.** Differences in abundance of genus adjusted as in ii.

417

418 **Figure 2.** *Escherichia/Shigella* variance stabilised abundances in twins discordant for rural-
419 urban living. Dashed lines link pairs.

420

421 **Supplementary figure legends**

422 **Figure S1.** Correlation of percentage greenspace within 800m, 3000m and 5000m of an
423 individual's residence.

424

425 **Figure S2.** Estimate greenspace percentage within 800m, 3000m and 5000m of an
426 individual's residence.

427

429 **Declarations**

430 **Ethics approval and consent to participate**

431 Favourable ethical opinion was granted by the formerly known St. Thomas' Hospital
432 Research Ethics Committee (REC). Following restructure and merging of REC, subsequent
433 amendments were approved by the NRES Committee London—Westminster (TwinsUK,
434 REC ref: EC04/015, 1 November 2011); use of microbiota samples was granted NRES
435 Committee London—Westminster (The Flora Twin Study, REC ref: 12/LO/0227, 1
436 November 2011).

437 **Availability of data and material**

438 The environmental exposures used in this paper were calculated thanks to several open-
439 source geographic datasets, as follows:

- 440 • The Rural-Urban classification:
 - 441 ○ England/Wales 2011 (Office for National Statistics):
442 [https://www.ons.gov.uk/methodology/geography/geographicalproducts/ruralur-](https://www.ons.gov.uk/methodology/geography/geographicalproducts/ruralurbanclassifications/2011ruralurbanclassification)
443 [banclassifications/2011ruralurbanclassification](https://www.ons.gov.uk/methodology/geography/geographicalproducts/ruralurbanclassifications/2011ruralurbanclassification)
 - 444 ○ Scotland 2016 (Scottish Government):
445 [https://www.gov.scot/publications/scottish-government-urban-rural-](https://www.gov.scot/publications/scottish-government-urban-rural-classification-2016/)
446 [classification-2016/](https://www.gov.scot/publications/scottish-government-urban-rural-classification-2016/)
- 447 • OS open roads 2017 (Ordnance Survey): [https://www.ordnancesurvey.co.uk/business-](https://www.ordnancesurvey.co.uk/business-government/products/open-map-roads)
448 [government/products/open-map-roads](https://www.ordnancesurvey.co.uk/business-government/products/open-map-roads)
- 449 • 25m² Land Cover Map of Great Britain 2015 (UK Centre for Ecology & Hydrology):
450 <https://www.ceh.ac.uk/services/land-cover-map-2015>

- 451 • The Index of Multiple Deprivation
- 452 ○ English (2019): [https://www.gov.uk/government/statistics/english-indices-of-](https://www.gov.uk/government/statistics/english-indices-of-deprivation)
- 453 deprivation
- 454 ○ Scottish (2016): <https://www2.gov.scot/Topics/Statistics/SIMD>
- 455 ○ Welsh (2019): [https://statswales.gov.wales/Catalogue/Community-Safety-](https://statswales.gov.wales/Catalogue/Community-Safety-and-Social70>Inclusion/Welsh-Index-of-Multiple-Deprivation/WIMD-2019)
- 456 and-Social70 Inclusion/Welsh-Index-of-Multiple-Deprivation/WIMD-2019

457 The European Bioinformatics Institute (EBI) accession numbers for the microbial DNA

458 sequences reported in this paper is ERP015317.

459 The processed ASV data, along with phenotypic data that can be used to recreate this analysis

460 is available following reasonable request to the TwinsUK data access committee. Information

461 on data access and how to apply is available at [http://www.twinsuk.ac.uk/data-](http://www.twinsuk.ac.uk/data-access/submission-procedure-2/)

462 [access/submission-procedure-2/](http://www.twinsuk.ac.uk/data-access/submission-procedure-2/). Please contact the corresponding author for further detail.

463 **Funding**

464 TwinsUK receives funding from the Wellcome Trust (WT081878MA), the National Institute

465 for Health Research (NIHR) Clinical Research Facility at Guy's & St Thomas' NHS

466 Foundation Trust and NIHR Biomedical Research Centre based at Guy's and St Thomas'

467 NHS Foundation Trust and King's College London. This work was also supported by the

468 Chronic Disease Research Foundation.

469 **Acknowledgements**

470 The authors of this paper wish to express our appreciation to all study participants of the

471 TwinsUK cohort for donating their samples and time. TwinsUK is funded by the Wellcome

472 Trust, Medical Research Council, European Union, the National Institute for Health Research

473 (NIHR)-funded BioResource, Clinical Research Facility and Biomedical Research Centre

474 based at Guy's and St Thomas' NHS Foundation Trust in partnership with King's College
475 London.

476 We thank Dr. Julia K Goodrich, Dr. Ruth E Ley and the Cornell technical team for generating
477 the microbial data..

478

479

480

481 References

482

483 Anderson, M.J., 2001. A new method for non-parametric multivariate analysis of variance. *Austral Ecol.* 26, 32–46.

484 <https://doi.org/10.1111/j.1442-9993.2001.01070.pp.x>

485 Belzer, C., Chia, L.W., Aalvink, S., Chamlagain, B., Piironen, V., Knol, J., de Vos, W.M., 2017. Microbial Metabolic

486 Networks at the Mucus Layer Lead to Diet-Independent Butyrate and Vitamin B12 Production by Intestinal

487 Symbionts. *MBio* 8, e00770-17. <https://doi.org/10.1128/mBio.00770-17>

488 Bowyer, R.C.E., Jackson, M.A., Le Roy, C.I., Lochlainn, M.N., Spector, T.D., Dowd, J.B., Steves, C.J., 2019.

489 Socioeconomic status and the gut microbiome: A twinsuk cohort study. *Microorganisms* 7.

490 <https://doi.org/10.3390/microorganisms7010017>

491 Bowyer, R.C.E., Jackson, M.A., Pallister, T., Skinner, J., Spector, T.D., Welch, A.A., Steves, C.J., 2018. Use of dietary

492 indices to control for diet in human gut microbiota studies. *Microbiome* 6. <https://doi.org/10.1186/s40168-018-0455-y>

493 Callahan, Benjamin J, McMurdie, P.J., Rosen, M.J., Han, A.W., Johnson, A.J.A., Holmes, S.P., 2016. DADA2: High-

494 resolution sample inference from Illumina amplicon data. *Nat. Methods* 13, 581–3.

495 <https://doi.org/10.1038/nmeth.3869>

496 Callahan, Ben J, Sankaran, K., Fukuyama, J.A., McMurdie, P.J., Holmes, S.P., 2016. Bioconductor Workflow for

497 Microbiome Data Analysis: from raw reads to community analyses. *F1000Research* 5, 1492.

498 <https://doi.org/10.12688/f1000research.8986.2>

499 Caporaso, J.G., Kuczynski, J., Stombaugh, J., Bittinger, K., Bushman, F.D., Costello, E.K., Fierer, N., Peña, A.G., Goodrich,

500 J.K., Gordon, J.I., Huttley, G.A., Kelley, S.T., Knights, D., Koenig, J.E., Ley, R.E., Lozupone, C.A., McDonald, D.,

501 Muegge, B.D., Pirrung, M., Reeder, J., Sevinsky, J.R., Turnbaugh, P.J., Walters, W.A., Widmann, J., Yatsunenko, T.,

502 Zaneveld, J., Knight, R., 2010. QIIME allows analysis of high-throughput community sequencing data. *Nat. Methods*

503 7, 335–336. <https://doi.org/10.1038/nmeth.f.303>

504 Claus, S.P., Guillou, H., Ellero-Simatos, S., 2016. The gut microbiota: A major player in the toxicity of environmental

505 pollutants? *npj Biofilms Microbiomes*. <https://doi.org/10.1038/npjbiofilms.2016.3>

506 Cresci, G.A., Bawden, E., 2015. Gut Microbiome: What We Do and Don't Know . *Nutr. Clin. Pract.* 30, 734–746.

507 <https://doi.org/10.1177/0884533615609899>

508 Dehingia, M., Thangjam devi, K., Talukdar, N.C., Talukdar, R., Reddy, N., Mande, S.S., Deka, M., Khan, M.R., 2016. Gut
509 bacterial diversity of the tribes of India and comparison with the worldwide data. *Sci. Rep.* 5, 18563.
510 <https://doi.org/10.1038/srep18563>

511 Falony, G., Joossens, M., Vieira-Silva, S., Wang, J., Darzi, Y., Faust, K., Kurilshikov, A., Bonder, M.J., Valles-Colomer,
512 M., Vandeputte, D., Tito, R.Y., Chaffron, S., Rymenans, L., Verspecht, C., De Sutter, L., Lima-Mendez, G., D'hoë,
513 K., Jonckheere, K., Homola, D., Garcia, R., Tigchelaar, E.F., Eeckhaut, L., Fu, J., Henckaerts, L., Zhernakova, A.,
514 Wijmenga, C., Raes, J., 2016. Population-level analysis of gut microbiome variation. *Science* (80-.). 352, 560–564.
515 <https://doi.org/10.1126/science.aad3503>

516 Geva-Zatorsky, N., Sefik, E., Kua, L., Pasman, L., Tan, T.G., Ortiz-Lopez, A., Yanortsang, T.B., Yang, L., Jupp, R., Mathis,
517 D., Benoist, C., Kasper, D.L., 2017. Mining the Human Gut Microbiota for Immunomodulatory Organisms. *Cell* 168,
518 928-943.e11. <https://doi.org/10.1016/J.CELL.2017.01.022>

519 Gloor, G.B., Macklaim, J.M., Pawlowsky-Glahn, V., Egozcue, J.J., 2017. Microbiome Datasets Are Compositional: And
520 This Is Not Optional. *Front. Microbiol.* 8, 2224. <https://doi.org/10.3389/fmicb.2017.02224>

521 Goodrich, J.K., Davenport, E.R., Beaumont, M., Jackson, M.A., Knight, R., Ober, C., Spector, T.D., Bell, J.T., Clark, A.G.,
522 Ley, R.E., 2016. Genetic Determinants of the Gut Microbiome in UK Twins. *Cell Host Microbe* 19, 731–743.
523 <https://doi.org/10.1016/j.chom.2016.04.017>

524 He, Y., Wu, W., Zheng, H.-M., Li, P., McDonald, D., Sheng, H.-F., Chen, M.-X., Chen, Z.-H., Ji, G.-Y., Zheng, Z.-D.-X.,
525 Mujagond, P., Chen, X.-J., Rong, Z.-H., Chen, P., Lyu, L.-Y., Wang, X., Wu, C.-B., Yu, N., Xu, Y.-J., Yin, J., Raes,
526 J., Knight, R., Ma, W.-J., Zhou, H.-W., 2018. Regional variation limits applications of healthy gut microbiome
527 reference ranges and disease models. *Nat. Med.* 24, 1532–1535. <https://doi.org/10.1038/s41591-018-0164-x>

528 Hillsdon, M., Coombes, E., Griew, P., Jones, A., 2015. An assessment of the relevance of the home neighbourhood for
529 understanding environmental influences on physical activity: how far from home do people roam? *Int. J. Behav. Nutr.*
530 *Phys. Act.* 12, 100. <https://doi.org/10.1186/s12966-015-0260-y>

531 Jackson, M.A., Goodrich, J.K., Maxan, M.-E., Freedberg, D.E., Abrams, J.A., Poole, A.C., Sutter, J.L., Welter, D., Ley,
532 R.E., Bell, J.T., Spector, T.D., Steves, C.J., 2016. Proton pump inhibitors alter the composition of the gut microbiota.
533 *Gut* 65, 749–56. <https://doi.org/10.1136/gutjnl-2015-310861>

534 Litvak, Y., Bäumlér, A.J., 2019. The founder hypothesis: A basis for microbiota resistance, diversity in taxa carriage, and
535 colonization resistance against pathogens. *PLOS Pathog.* 15, e1007563. <https://doi.org/10.1371/journal.ppat.1007563>

536 Love, M.I., Huber, W., Anders, S., 2014. Moderated estimation of fold change and dispersion for RNA-seq data with

- 537 DESeq2. *Genome Biol.* 15, 550. <https://doi.org/10.1186/s13059-014-0550-8>
- 538 Magnúsdóttir, S., Ravcheev, D., de Crécy-Lagard, V., Thiele, I., 2015. Systematic genome assessment of B-vitamin
539 biosynthesis suggests co-operation among gut microbes. *Front. Genet.* 6, 148.
540 <https://doi.org/10.3389/fgene.2015.00148>
- 541 Marchesi, J.R., Ravel, J., 2015. The vocabulary of microbiome research: a proposal. *Microbiome* 3, 31.
542 <https://doi.org/10.1186/s40168-015-0094-5>
- 543 McMurdie, P.J., Holmes, S., Kindt, R., Legendre, P., O'Hara, R., 2013. phyloseq: An R Package for Reproducible
544 Interactive Analysis and Graphics of Microbiome Census Data. *PLoS One* 8, e61217.
545 <https://doi.org/10.1371/journal.pone.0061217>
- 546 Obregon-Tito, A.J., Tito, R.Y., Metcalf, J., Sankaranarayanan, K., Clemente, J.C., Ursell, L.K., Zech Xu, Z., Van Treuren,
547 W., Knight, R., Gaffney, P.M., Spicer, P., Lawson, P., Marin-Reyes, L., Trujillo-Villarreal, O., Foster, M., Gujja-
548 Poma, E., Troncoso-Corzo, L., Warinner, C., Ozga, A.T., Lewis, C.M., 2015. Subsistence strategies in traditional
549 societies distinguish gut microbiomes. *Nat. Commun.* 6, 6505. <https://doi.org/10.1038/ncomms7505>
- 550 Oksanen, J., Blanchet, F.G., Kindt, R., Legendre, P., Minchin, P.R., O'Hara, R.B., Simpson, G.L., Solymos, P., Stevens,
551 M.H.H., Wagner, H., 2015. vegan: Community Ecology Package. R package version 2.0-10. URL: [https://CRAN.R-](https://CRAN.R-project.org/package=vegan)
552 [project.org/package=vegan](https://CRAN.R-project.org/package=vegan)
- 553 Ordiz, M.I., May, T.D., Mihindikulasuriya, K., Martin, J., Crowley, J., Tarr, P.I., Ryan, K., Mortimer, E., Gopalsamy, G.,
554 Maleta, K., Mitreva, M., Young, G., Manary, M.J., 2015. The effect of dietary resistant starch type 2 on the
555 microbiota and markers of gut inflammation in rural Malawi children. *Microbiome* 3, 37.
556 <https://doi.org/10.1186/s40168-015-0102-9>
- 557 Palarea-Albaladejo, J., Martín-Fernández, J.A., 2015. zCompositions — R package for multivariate imputation of left-
558 censored data under a compositional approach. *Chemom. Intell. Lab. Syst.* 143, 85–96.
559 <https://doi.org/10.1016/j.chemolab.2015.02.019>
- 560 Robinson, J.M., Jorgensen, A., 2020. Rekindling old friendships in new landscapes: The environment–microbiome–health
561 axis in the realms of landscape research. *People Nat.* 2, 339–349. <https://doi.org/10.1002/pan3.10082>
- 562 Rook, G.A., 2013. Regulation of the immune system by biodiversity from the natural environment: an ecosystem service
563 essential to health. *Proc. Natl. Acad. Sci. U. S. A.* 110, 18360–18367. <https://doi.org/10.1073/pnas.1313731110>
- 564 Rook, G.A.W., Lowry, C.A., Raison, C.L., 2013. Microbial “Old Friends”, immunoregulation and stress resilience.

565 Evolution (N. Y). 46–64. <https://doi.org/10.1093/emph/eot004>

566 Rothschild, D., Weissbrod, O., Barkan, E., Kurilshikov, A., Korem, T., Zeevi, D., Costea, P.I., Godneva, A., Kalka, I.N.,
567 Bar, N., Shilo, S., Lador, D., Vila, A.V., Zmora, N., Pevsner-Fischer, M., Israeli, D., Kosower, N., Malka, G., Wolf,
568 B.C., Avnit-Sagi, T., Lotan-Pompan, M., Weinberger, A., Halpern, Z., Carmi, S., Fu, J., Wijmenga, C., Zhernakova,
569 A., Elinav, E., Segal, E., 2018. Environment dominates over host genetics in shaping human gut microbiota. *Nature*
570 555, 210–215. <https://doi.org/10.1038/nature25973>

571 Schmidt, T.S.B., Raes, J., Bork, P., 2018. The Human Gut Microbiome: From Association to Modulation. *Cell* 172, 1198–
572 1215. <https://doi.org/10.1016/j.cell.2018.02.044>

573 Searle, S.D., Mitnitski, A., Gahbauer, E.A., Gill, T.M., Rockwood, K., 2008. A standard procedure for creating a frailty
574 index. *BMC Geriatr.* 8, 1–10. <https://doi.org/10.1186/1471-2318-8-24>

575 Shreiner, A.B., Kao, J.Y., Young, V.B., 2015. The gut microbiome in health and in disease. *Curr. Opin. Gastroenterol.* 31,
576 69–75. <https://doi.org/10.1097/MOG.0000000000000139>

577 Stamper, C.E., Hoisington, A.J., Gomez, O.M., Halweg-Edwards, A.L., Smith, D.G., Bates, K.L., Kinney, K.A., Postolache,
578 T.T., Brenner, L.A., Rook, G.A.W., Lowry, C.A., 2016. The Microbiome of the Built Environment and Human
579 Behavior: Implications for Emotional Health and Well-Being in Postmodern Western Societies. *Int. Rev. Neurobiol.*
580 131, 289–323. <https://doi.org/10.1016/BS.IRN.2016.07.006>

581 Stekhoven, D.J., Buhlmann, P., 2012. MissForest--non-parametric missing value imputation for mixed-type data.
582 *Bioinformatics* 28, 112–118. <https://doi.org/10.1093/bioinformatics/btr597>

583 Tringe, S.G., Hugenholtz, P., 2008. A renaissance for the pioneering 16S rRNA gene. *Curr. Opin. Microbiol.* 11, 442–446.
584 <https://doi.org/10.1016/J.MIB.2008.09.011>

585 Twohig-Bennett, C., Jones, A., 2018. The health benefits of the great outdoors: A systematic review and meta-analysis of
586 greenspace exposure and health outcomes. *Environ. Res.* 166, 628–637. <https://doi.org/10.1016/j.envres.2018.06.030>

587 Tyakht, A. V., Alexeev, D.G., Popenko, A.S., Kostyukova, E.S., Govorun, V.M., 2014. Rural and urban microbiota. *Gut*
588 *Microbes* 5, 351–356. <https://doi.org/10.4161/gmic.28685>

589 Eryk Walczak (2021). PostcodesioR: API Wrapper Around 'Postcodes.io'. R package version 0.3.0. URL: [https://CRAN.R-](https://CRAN.R-project.org/package=PostcodesioR)
590 [project.org/package=PostcodesioR](https://CRAN.R-project.org/package=PostcodesioR).

591 Weersma, R.K., Zhernakova, A., Fu, J., 2020. Interaction between drugs and the gut microbiome. *Gut* 69, 1510 LP – 1519.

592 <https://doi.org/10.1136/gutjnl-2019-320204>

593

594