

1 A quantitative evaluation of MIRU-VNTR typing against whole-genome sequencing for identifying  
2 *Mycobacterium tuberculosis* transmission: A prospective observational cohort study

3 David Wyllie<sup>1,2,3\*</sup>, Jennifer Davidson<sup>4</sup>, Tim Walker<sup>1</sup>, Preeti Rathod<sup>5</sup>, Derrick Crook<sup>1,3</sup>, Tim Peto<sup>1,3</sup>,  
4 Esther Robinson<sup>5</sup>, Grace Smith<sup>5</sup>, Colin Campbell<sup>4</sup>

5  
6 **Affiliations**

7 <sup>1</sup> Nuffield Department of Medicine, John Radcliffe Hospital, Headley Way, Oxford OX3 9DU, UK

8 <sup>2</sup> Public Health England Academic Collaborating Centre, John Radcliffe Hospital, Headley Way,  
9 Oxford OX3 9DU, UK

10 <sup>3</sup> The National Institute for Health Research Health Protection Research Unit (NIHR HPRU) in  
11 Healthcare Associated Infections and Antimicrobial Resistance at University of Oxford

12 <sup>4</sup> Tuberculosis Section, National Infection Service, Public Health England, 61 Colindale Avenue,  
13 London NW9 5EQ, UK

14 <sup>5</sup> Public Health England National Regional Mycobacteriology Laboratory North and Midlands,  
15 Heartlands Hospital, Birmingham BS9 5SS

16

17 \* correspondence

18

19 Summary

20 *Background*

21 Mycobacterial Interspersed Repetitive Unit-Variable Number Tandem Repeat (MIRU-VNTR) typing is  
22 widely used in high-income countries for *Mycobacterium tuberculosis* typing. Whole-genome  
23 sequencing (WGS) is known to deliver greater specificity, but no quantitative prospective comparison  
24 has yet been undertaken.

25 *Methods*

26 We studied isolates from the English Midlands, sampled consecutively between 1 January 2012 and  
27 31 December 2015. In addition to routinely performed MIRU-VNTR typing, DNA was extracted from  
28 liquid cultures and sequenced using Illumina technology. Demographic and epidemiological data were  
29 extracted from the Enhanced Tuberculosis Surveillance system maintained by Public Health England.  
30 Closely related samples, defined using a threshold of five single nucleotide variants (SNVs), were  
31 compared to samples with identical MIRU-VNTR profiles, with shared epidemiological risk factors,  
32 and to those with both characteristics.

33 *Findings*

34 1,999 patients were identified for whom at least one *M. tuberculosis* isolate had been MIRU-VNTR  
35 typed and sequenced. Comparing epidemiological risk factors with close genetic relatedness, only co-  
36 residence had a positive predictive value of over 5%. Excluding co-resident individuals, 18.6% of  
37 patients with identical MIRU-VNTR profiles were within 5 SNVs. Where patients also shared social  
38 risk factors and ethnic group, this rose to 48%. Only 8% of MIRU-VNTR linked pairs in lineage 1 were  
39 within 5 SNV, compared to 31% in lineage 4.

40 *Interpretation*

41 In the setting studied, MIRU-VNTR typing and epidemiological risk factors are poorly predictive of  
42 close genomic relatedness, assessed by SNV. MIRU-VNTR performance varies markedly by lineage.

43 *Funding*

44 Public Health England, National Institute of Health Research Oxford Biomedical Research Centre.

45

46

47

48 Research in context

49 *Evidence before this study*

50 We searched Pubmed using the search terms ‘whole genome sequencing’ and ‘MIRU-VNTR’ and  
51 ‘tuberculosis’ for English language articles published up to December 21<sup>st</sup>, 2017. Multiple studies  
52 have shown that most pairwise genomic comparisons will be within five SNVs when direct  
53 transmission has occurred from one individual to another. Both outbreak studies and population  
54 studies have demonstrated how whole-genome sequencing generates smaller clusters than MIRU-  
55 VNTR typing, and how sequence data allows for differentiation of isolates within a cluster. However,  
56 no systematic comparison of MIRU-VNTR typing vs. WGS has however been published. The degree  
57 to which WGS provides more specific results, and the degree to which it is likely to be more cost  
58 effective, therefore remains uncertain.

59 *Added value of this study*

60 This study seeks to quantify the predictive value of identical MIRU-VNTR profiles, and of overlapping  
61 demographic and epidemiological data, for close genomic relatedness in a cosmopolitan setting.  
62 Importantly, it demonstrates that in our setting MIRU-VNTR-based clustering predicts genomic  
63 relatedness differently depending on *M. tuberculosis* lineage. Whether this is due to biological  
64 differences between the lineages or to immigration patterns, it is likely that these findings are relevant  
65 to other cosmopolitan settings. These data provide an explanation as to why MIRU-VNTR typing was  
66 not cost-effective when implemented in England, and indicate WGS may perform substantially better.

67 *Implications of all the available evidence*

68 Whilst it is generally accepted that WGS provides more informative results than MIRU-VNTR typing,  
69 the latter is still practiced widely under the belief that it remains a helpful tool for public health  
70 investigations. This study shows that whilst differing MIRU-VNTR profiles help exclude close genomic  
71 relatedness, matching profiles rarely predict such relatedness. Having quantified its predictive value at  
72 a population level, this study should hasten the transition from MIRU-VNTR typing to WGS in other  
73 settings similar to ours.

74

75 Introduction

76 In 2016 there were 5,664 notified cases of tuberculosis in the England, with an incidence of 10.2 per  
77 100,000 population.<sup>1</sup> Despite a steady fall in incidence since its peak early this decade, this remains  
78 the highest rate in western Europe, outside of the Iberian peninsula.<sup>2</sup> Much of the recent decline in  
79 incidence has been due to a falling number of patients born outside of the UK. However, this decline  
80 slowed in 2016, with domestic transmission likely to still be contributing towards the residual case  
81 load.

82 Rapid detection of *Mycobacterium tuberculosis* transmission should offer enhanced opportunities for  
83 disease control.<sup>3,4</sup> In England, as in many high-income countries, tuberculosis transmission has been  
84 identified with the help of Mycobacterial Interspersed Repetitive Unit-Variable Number Tandem  
85 Repeat (MIRU-VNTR) typing, which clusters cultured isolates on the basis of their molecular  
86 fingerprints.<sup>5,6</sup> A recent post-deployment evaluation of the MIRU-VNTR-based surveillance  
87 programme in England has however questioned the cost-effectiveness of this approach.<sup>7</sup>

88 Since 2015, Public Health England has been undertaking a phased introduction of routine whole  
89 genome sequencing (WGS) for all mycobacterial cultures.<sup>8</sup> This has meant the relatedness of isolates  
90 could be simultaneously compared using both single nucleotide variants (SNV) and by MIRU-VNTR  
91 typing, and has provided a novel opportunity to compare the added value of whole genome  
92 sequencing (<sup>9-14</sup>;Table 1) in an unselected population, at scale.

93 Here we estimate what proportion of *M. tuberculosis* isolates from a cosmopolitan area of central  
94 England that are linked by MIRU-VNTR typing, or have associated epidemiological risk factors, are  
95 closely genomically related.

96

97

98 Methods

99 *Samples studied for comparison of MIRU-VNTR with SNVs*

100 Consecutive *M. tuberculosis* isolates from the Public Health England Centre for Regional  
101 Mycobacteriology Laboratory, Birmingham between 1 January 2012 and 31 December 2015 were  
102 included in the study. This laboratory serves a large catchment of approximately 12 million persons in  
103 the English Midlands, a region which includes high, medium (40-150 cases per 100,000 population),  
104 and low TB incidence areas.

105 *Identification and MIRU-VNTR typing*

106 Clinical samples were grown in Mycobacterial Growth Indicator tubes (MGIT) (Becton Dickinson, New  
107 Jersey, USA), and *M. tuberculosis* was identified using Ziehl-Neelsen staining, followed by nucleic  
108 acid amplification and hybridisation using Genotype Mycobacterium CM hybridisation tests (Hain  
109 LifeScience, Nehren, Germany). MIRU-VNTR typing<sup>5</sup> was performed on the first isolate from each  
110 patient in each calendar year, following protocols then in place.

111 *Laboratory and bioinformatic processing*

112 This was carried out as described.<sup>10</sup> Nucleic acid was extracted from 1-7 ml of MGIT culture as  
113 described.<sup>8</sup> Illumina 150 bp paired end DNA libraries were made using Nextera XT version 2  
114 chemistry kits and sequenced on MiSeq instruments (Illumina). Reads were mapped to the H37Rv v2  
115 reference genome (Genbank: NC000962.2) using Stampy<sup>15</sup>, and aligned to Bam files parsed with  
116 Samtools mPileup<sup>16</sup>, with further filtering performed based on the base and alignment quality (q30 and  
117 Q30 cutoffs, respectively). SNV variation was reported but indels were not considered as part of this  
118 work as they have been reported to be less reliably called than SNVs.<sup>15</sup> Bases supported only by low  
119 confidence base calls were recorded as uncertain ('N'), as were positions with > 10% minor variant  
120 frequencies, and all calls at the genomic positions included in Supplementary Data 1, since these  
121 regions were repetitive (as identified by self-self blastn analysis) or were found to commonly contain  
122 low-confidence mapping (*rrl*, *rrs*, *rpoC* and *Rv2082* loci). Such uncertain bases were ignored in  
123 pairwise SNV computations.

124 *Metrics of relatedness*

125 We used pairwise SNV distances between isolates as a metric of close genetic relatedness,  
126 considering isolates closely genetically related when their pairwise SNV distance was less a particular  
127 SNV threshold. For the main analysis, 5 SNV was used as the threshold, but a range of other  
128 thresholds were considered in sensitivity analyses.

129 Lineage assignment was performed using ancestral SNVs, as described.<sup>17</sup> Relatedness between  
130 samples was determined by comparing the number of mismatching positions between loci using  
131 BugMat.<sup>18</sup> Relatedness between MIRU-VNTR profiles compared the total number of differences in  
132 repeat lengths at each of the 24 loci. For example, for a one-locus typing scheme, if isolate 1 had 3  
133 repeats, and isolate 2 had 5 repeats, we coded this as a 2 MIRU-VNTR repeat unit difference.

134 *Collection and collation of patient data*

135 Demographic data (sex, age, ethnic group and residence), and social risk factor data (current or  
136 history of imprisonment, drug misuse, alcohol misuse or homelessness) were obtained from the  
137 Enhanced Tuberculosis Surveillance system. Co-residence was defined as having the same first line  
138 of address and postcode.

139 *Statistical analyses*

140 We considered a series of categorical variables as predictors of close genomic relatedness in logistic  
141 regression analyses. Additionally, for some variables, we constructed composite categorical variables  
142 reflecting whether more than one risk factor was present. For each given SNV threshold, we  
143 estimated odds ratios for close genomic relatedness using logistic regression. Separately, we  
144 modelled the relationship between SNV variation (s) (outcome), *Mycobacterium tuberculosis* lineage  
145 (l, a discrete variable) and n, the number of MIRU-VNTR repeat number differences observed, as  
146 defined above. We modelled

147  $E(s) \sim n + l + n * l$

148 thus allowing estimation of both lineage-specific variation in the absence of any variation in MIRU-  
149 VNTR types, and how SNV increased with increasing MIRU-VNTR differences. We used quantile  
150 regression (R quantreg package) for the main analysis as homoscedascity assumptions were  
151 violated. All analyses used R 3.3.1 for Windows.

152 *Ethical framework*

153 Public health action taken as a result of notification and surveillance is one of the Public Health  
154 England's key roles as stated in the Health and Social Care Act 2012 and subsequent Government  
155 directives which provide the mandate and legislative basis to undertake necessary follow-up. Part of  
156 this follow-up is identification of epidemiological and molecular links between cases. This work is  
157 part of service development carried out under this framework, and as such explicit ethical approval is  
158 unnecessary.

159 *Funding source*

160 This study is supported by the Health Innovation Challenge Fund (a parallel funding partnership  
161 between the Wellcome Trust [WT098615/Z/12/Z] and the Department of Health [grant HICF-T5-358])  
162 and NIHR Oxford Biomedical Research Centre. Professor Derrick Crook is affiliated to the National  
163 Institute for Health Research Health Protection Research Unit (NIHR HPRU) in Healthcare Associated  
164 Infections and Antimicrobial Resistance at University of Oxford in partnership with Public Health  
165 England. Professor Crook is based at University of Oxford. The views expressed are those of the  
166 author(s) and not necessarily those of the NHS, the NIHR, the Department of Health or Public Health  
167 England. The sponsors of the study had no role in study design, data collection, data analysis, data  
168 interpretation, or writing of the report. The corresponding author had full access to all the data in the  
169 study and had final responsibility for the decision to submit for publication.

170 Results

171 *Isolates studied*

172 We studied all *M. tuberculosis* isolates consecutively grown in, or referred to, the Public Health  
173 England Mycobacterial reference centre for the English Midlands between 2012-2015 (n=2,718)  
174 (Figure 1). We excluded 551 isolates because MIRU-VNTR typing had already been performed on  
175 another isolate (protocol was to MIRU-VNTR type one isolate per patient per year), and 57 isolates  
176 because of technical concerns about laboratory processing (Figure 1). The remaining 2,110 isolates  
177 came from 2,020 discrete patients. A further 16 isolates were excluded because multiple isolates from  
178 the same individual were separated by >12 single nucleotide variants (SNVs) (suggestive of technical  
179 error), along with five recurrent cases of *M. tuberculosis* infection, leaving 1,999 isolates each derived  
180 from a different patient.

181 There were more male than female patients (1176, 58%). 1155 (58%) were aged between 15-44  
182 years old. 1325 patients (66%) were born outside the UK and 1437 (71%) were of non-White ethnicity  
183 (Table 2). *M. tuberculosis* lineage 4 (Euro-American) was the most commonly isolated lineage  
184 (n=954, 48%) with lineages 1, 2, and 3 also commonly represented (176 (9%), 137 (7%), 704 (35%)  
185 isolates respectively) (Table 2). *M. tuberculosis* lineage was associated with country of birth, with  
186 lineage 3 being most common in individuals born in India or Pakistan (Table 3).

187 *Epidemiological risk factors and the prediction of close relatedness*

188 Using pairwise SNV distances within 5 SNVs between isolates to define genomic relatedness, we  
189 determined how various shared epidemiological data altered the odds of relatedness. Figure 2A  
190 shows estimated odds ratios of close genomic relatedness, in the presence, relative to the absence,  
191 of a series of risk factors. The proportion of paired isolates that are closely genomically related, given  
192 a particular risk factor, was also calculated. This represents the positive predictive value (PPV) of  
193 each risk factor. SNV thresholds other than 5 SNVs were analysed in sensitivity analyses (Web extra  
194 Fig. S1-S6), with similar results.

195 Predictably, residence at the same address was most strongly associated with close genomic  
196 relatedness (OR 8,000, 95% CI 5,000, 13,000). This corresponds to a PPV of 42%, indicating the  
197 majority of co-resident cases in this series were not closely genomically related, something discussed  
198 below. However, it was rare for two patients to share an address, with only 85 isolates derived from  
199 such settings. Other risk factors studied included sharing a self-identified ethnic group with another  
200 patient or being in a similar age bracket. Both were weakly associated with genomic relatedness  
201 (estimated odds ratios of 10 or less), with the highest risk of close genomic relatedness for an ethnic  
202 group seen for the smallest ethnic group studied (those identifying as Black Caribbean or Black Other;  
203 n=71; OR 16, 95% CI 8, 32). Similarly, there was a modest increase in the odds of close genomic  
204 relatedness where two isolates were from individuals with social risk factors (current or history of  
205 imprisonment, drug misuse, alcohol misuse or homelessness) (OR 9, 95% CI 4, 16). In all these  
206 cases however, the PPV was less than 1%.

207 *MIRU-VNTR profiles as predictors of close relatedness*

208 Having identical MIRU-VNTR profiles conferred an odds ratio of close genomic relatedness of 2,800  
209 (95% CI 2,200, 3,400) on paired isolates, compared with paired isolates with different MIRU-VNTR  
210 profiles, with an associated 18.6% PPV (Figure 2B). With 1 locus discordant, the corresponding odds  
211 ratio and PPV were much lower (OR 210, 95% CI 160,270; PPV 1.7%).

212 To understand how MIRU-VNTR profile and epidemiological data can complement each other in the  
213 identification of close relatedness, we assessed combinations of the presence of identical MIRU-  
214 VNTR profiles, social risk factors, and shared ethnicity, all factors which are significantly associated  
215 with close relatedness individually (Figure 2A). Excluding individuals who were resident at the same  
216 address, identical MIRU-VNTR profile was more predictive of close relatedness when shared risk  
217 factors were present, but for all the combinations studied the PPV remained low (15%, 18%, 33%,  
218 48% with no shared risk factors, same ethnic group but no social risk factors, shared social risk  
219 factors but different ethnic group, and both shared ethnic group and social risk factors, respectively).

220 *SNV - MIRU-VNTR relationships vary by lineage*

221 While MIRU-VNTR profiles predict close genetic relatedness (defined by SNVs) better than most  
222 social risk factors (Figure 2A), we observed that the PPV differs markedly by *M. tuberculosis* lineage  
223 (Fig. 3). For lineages 1, 2, 3 and 4, which together account for 1,977/1,999 (99%) of the isolates  
224 studied, we compared pairwise comparisons within each lineage by MIRU-VNTR similarity (Fig. 4).  
225 For lineages 1 and 4, pairwise SNV distances increased over the range 0 to 8 MIRU-VNTR unit  
226 differences, until at higher MIRU-VNTR distances the pairwise distances approximated the within-  
227 lineage median pairwise SNV distance (Fig. 4). For lineages 2 and 3 the median was reached by 3  
228 MIRU-VNTR differences. Overall there was less variation between paired isolates within lineages 2  
229 and 3 (median pairwise distances 205 and 334, respectively) compared to paired isolates within  
230 lineages 1 and 4 (median pairwise distances 840, and 685). However, for paired isolates differing by  
231 between zero and 4 MIRU-VNTR loci, the least variation was seen within lineage 4.

232 To quantify how the relationship between MIRU-VNTR and SNVs differed by lineage, we modelled  
233 SNV distances between paired isolates, assuming a linear relationship with MIRU-VNTR profile  
234 distances over the range of 0-3 MIRU-VNTR unit differences (Figure 4, red dots show fitted medians,  
235 and Supplementary Data 2). For lineage 4 isolates, among pairs with identical MIRU-VNTR profiles,  
236 there was a median of  $10 \pm 0.4$  SNV (median  $\pm$  standard error). For paired isolates with identical  
237 MIRU-VNTR profiles in lineages 1, 2, and 3, SNV distances were  $122 \pm 21$ ,  $159 \pm 3$ , and  $82 \pm 3$   
238 (median  $\pm$  standard error), respectively. According to current estimates of *M. tuberculosis* clock rates,  
239 these correspond to about 250, 300, and 150 years of evolution, respectively, compared to about 20  
240 years for lineage 4<sup>19</sup>.

241 For each MIRU-VNTR unit difference in lineage 4, there was a median increase of  $59 \pm 0.6$  SNV. For  
242 lineage 1, a similar increase in SNV with increasing MIRU-VNTR differences was observed to that in  
243 lineage 4 (het.  $p = 0.32$ ), whereas for lineages 2 and 3 the relationship was very different from lineage  
244 4 (het.  $p < 10^{-20}$  for both comparisons). Indeed, for paired isolates in lineage 2, SNVs were only



245 weakly associated with MIRU-VNTR distance. Thus, in the population studied, the performance of  
246 MIRU-VNTR profiles in defining evolutionarily related groups differed between lineage 4 (Euro-  
247 American) isolates, and lineages 1, 2 and 3.

248

249

250 Discussion

251 In this prospective study of a cosmopolitan population in the English Midlands, we have quantified  
252 how well recent transmission, as defined a 5 SNV threshold, is predicted by shared epidemiological  
253 risk factors, by MIRU-VNTR typing, or by a combination of both.<sup>19</sup> We have also demonstrated how  
254 lineage strongly affects the performance of MIRU-VNTR-based predictions.

255 Overall, the PPV for recent transmission for any two isolates with an identical MIRU-VNTR type was  
256 only 18.6%. Excluding cases resident at the same address, the PPV varied from as low as 14.8% to  
257 48.0% if shared risk factors were present alongside identical MIRU-VNTR profiles (Figure 2).

258 However, PPVs for shared MIRU-VNTR profiles differed significantly by lineage, with the strongest  
259 associations seen in lineage 4 (European-American), which was also most frequently observed  
260 lineage in the Midlands. The number of patient-to-patient links that need to be investigated to find a  
261 single case of recent transmission between non-co-resident individuals with shared MIRU-VNTR  
262 types is thus between two and seven, depending on the presence of shared social risk factors.

263 These data demonstrate that the previous routine practice of grouping samples based on MIRU-  
264 VNTR identity, or on a combination of MIRU-VNTR identity and shared epidemiological risk factors,  
265 generates highly heterogeneous results, and is likely to contribute to the low cost-effectiveness of  
266 MIRU-VNTR typing.<sup>7</sup> Importantly, our data also demonstrate how lineage markedly affects the PPV of  
267 MIRU-VNTR links, with the best results seen for lineage 4. To our knowledge, lineage has not been  
268 routinely taken into consideration when matching isolates by MIRU-VNTR for surveillance reasons.

269 One possible explanation for why SNV distances between paired isolates sharing a MIRU-VNTR  
270 profile within lineages 1, 2 and 3 were greater than for lineage 4 is that the Indo-Oceanic, East-Asian  
271 (including Beijing) and East-African Indian lineages are more endemic to countries other than the UK,  
272 and that patients diagnosed with these tuberculosis lineages in the UK were infected overseas. Were  
273 this the case, closely genomically related strains would be less likely to be found in England. For  
274 example, lineage 3 isolates were most common in individuals born in India and Pakistan, relative to  
275 other individuals, supporting this hypothesis (Table 3). A second possible explanation is that the rate  
276 of diversification of MIRU-VNTR types relative to SNVs differs between major lineages. Thirdly,  
277 MIRU-VNTR variation can result in the same profile via different evolutionary routes (homoplasy)<sup>20</sup>, a  
278 phenomenon which could also explain the rather flat relationship observed between MIRU-VNTR  
279 distance and SNV distance seen in lineages 2 and 3. Whatever the mechanism(s) operating, our data  
280 implies that the lineages, and their epidemiology, may influence the wide variation in the proportion of  
281 TB cases clustering using MIRU-VNTR profiling reported in different settings.<sup>9,21</sup>

282 It was surprising to us that among individuals resident at the same address, only 42% of these pairs  
283 were closely genomically linked. One explanation for this relatively low proportion is that some  
284 patients from highly endemic countries are likely to co-habit with others from highly endemic  
285 countries, potentially increasing the chances of non-clustered isolates, originating from separate  
286 exposures, being linked to the same address. Another scenario that could lead to a similar effect  
287 would be UK born patients with multiple social risk factors sharing hostels. In both settings, co-

288 resident individuals with TB would be expected to have an increased risk of having acquired their  
289 infection from individuals in high prevalence populations with whom they have been in contact outside  
290 the residential setting.

291 One limitation to this study is that is that the results cannot necessarily be generalised to other  
292 settings with different patterns of transmission, rates of disease, patterns of immigration, and relative  
293 prevalence of different lineages. However, the region studied was large and included a mixture of  
294 incidence areas, and both urban and rural settings. Another potential limitation is that we cannot be  
295 sure that risk factor data was recorded in a fully sensitive manner. Under-ascertainment of risk factor  
296 data would reduce the apparent contribution of risk factor data to identifying close genetic neighbours.  
297 However, even in the population in which we found in which MIRU-VNTR profiling works best (lineage  
298 4 infections), and in subjects for whom shared risk factors were recorded, the combination of MIRU-  
299 VNTR identity and shared risk factors only detects about one in two closely related isolate pairs.

300 In summary, these data help quantify the limitations of MIRU-VNTR typing for tuberculosis  
301 transmission surveillance and control. With routine diagnostic services beginning to transition to WGS  
302 technology in multiple high incidence countries, as England already has, our data indicates one can  
303 expect to see a reduction in the number of potential links requiring epidemiological investigation by a  
304 factor of about five. WGS thus stands a much greater chance of contributing to a cost effective control  
305 program than MIRU-VNTR typing in low-burden, cosmopolitan settings such as ours.

306

307 Tables

308 *Table 1 Previous studies including both MIRU-VNTR and SNV analysis of M. tuberculosis*

Samples	Comment	Reference
36 archived Manila strain isolates	SNV analysis revealed variation not demonstrated by MIRU-VNTR.	<sup>9</sup>
390 retrospective isolates from the English Midlands	Genetic heterogeneity within MIRU-VNTR clusters demonstrated. 5 and 12 SNV proposed as potential cut offs for epidemiological relatedness.	<sup>10</sup>
199 epidemiologically linked cases sequenced retrospectively	Relationship with MIRU-VNTR profile was not addressed	<sup>22</sup>
36 isolates from an outbreak	SNV analysis revealed variation not demonstrated by MIRU-VNTR.	<sup>23</sup>
50 cases from an outbreak	SNV analysis revealed variation not demonstrated by MIRU-VNTR.	<sup>11</sup>
1,000 isolate sample of 2,248. Representative of Russian population studied, plus 28 diverse sequences	Relationship with MIRU-VNTR profile was not addressed. Multiple sub-lineages observed within Lineage 4 (Euro-American).	<sup>24</sup>
69 cases from an outbreak defined by a SNV	SNV analysis revealed variation not demonstrated by MIRU-VNTR.	<sup>12</sup>
86 cases from an outbreak	SNV analysis revealed variation not demonstrated by MIRU-VNTR.	<sup>13</sup>
90 cases belonging to 35 MIRU-VNTR clusters	MIRU-VNTR performance overestimated transmission particularly in immigrants infected with closely related strains	<sup>14</sup>

309

310

311 *Table 2 Details of Samples studied*

Category	Property	Number of samples
Number of social risk factors (homelessness, prison, alcohol use, drug use)	0	1761
	1	136
	2	51
	3	26
	4	3
	Not available	22
Gender	Female	801
	Male	1176
	Not available	22
Age group	0-14	45
	15-44	1155
	45-64	442
	65+	335
	Not available	22
Year sample taken	2007	1
	2010	1
	2011	5
	2012	355
	2013	584
	2014	507
	2015	524
	Not available	22
PHE Region of patient's residence	London	6
	Midlands & East of England	1721
	North of England	243
	South of England	3
	Not available	26
Self-declared ethnic group	Bangladeshi	31
	Black-African	267
	Black-Caribbean	57
	Black-Other	14
	Chinese	29
	Indian	564
	Mixed / Other	143
	Pakistani	332
	White	508
	Not available	54
UK Born	Non-UK Born	1325
	UK Born	592
	Not available	82

312

313 *Table 3 Lineage of isolates studied*

Place of birth	Lineage					Total
	1	2	3	4	Other	
UNITED KINGDOM	19 (3.2%)	33 (5.5%)	136 (23%)	391 (66%)	13 (2.1%)	592 (100%)
INDIA	81 (18%)	18 (4.0%)	246 (55%)	102 (23%)	1 (0.2%)	448 (100%)
PAKISTAN	16 (6.3%)	6 (2.3%)	178 (71%)	51 (20%)	1 (0.4%)	252 (100%)
SOMALIA	8 (15%)	2 (3.8%)	24 (45%)	18 (33%)	1 (1.9%)	53 (100%)
ZIMBABWE	3 (6.0%)	7 (14%)	1 (2.0%)	37 (76%)	1 (2.0%)	49 (100%)
ERITREA	3 (6.5%)	2 (4.3%)	16 (35%)	25 (54%)	0 (0.0%)	46 (100%)
POLAND	0 (0.0%)	1 (2.8%)	1 (2.8%)	34 (94%)	0 (0.0%)	36 (100%)
ROMANIA	0 (0.0%)	0 (0.0%)	0 (0.0%)	28 (100%)	0 (0.0%)	28 (100%)
LITHUANIA	0 (0.0%)	8 (33%)	0 (0.0%)	16 (66%)	0 (0.0%)	24 (100%)
Other	36 (9.5%)	58 (15%)	66 (18%)	208 (55%)	9 (2.4%)	377 (100%)
Not known	10 (10%)	2 (2.1%)	36 (38%)	44 (49%)	2 (2.1%)	94 (100%)
<b>Total</b>	<b>176</b>	<b>137</b>	<b>704</b>	<b>954</b>	<b>28</b>	<b>1999</b>

314

315 Table Legends

316 *Table 1*

317 Published comparisons between MIRU-VNTR and SNV based M. tuberculosis typing.

318 *Table 2*

319 Details of isolates studied.

320 *Table 3*

321 Lineages of isolates studied, and recorded country of birth of the subjects.

322

323 Figure Legends

324 *Figure 1 Flowchart showing the samples studied*

325 Flowchart showing the samples studied.

326 *Figure 2 Relationship between MIRU-VNTR profile, epidemiological risk factors and genetic*  
327 *relatedness*

328 The odds ratio predicting closely related isolates (defined by having five or fewer single nucleotide  
329 variants between them) associated with sharing a series of epidemiological properties. PPV denotes  
330 positive predictive values. n refers to the number of subjects having the property described. For  
331 example, there were 801 female subjects.

332 *Figure 3 Association between lineage, close genetic relatedness and MIRU-VNTR profile*

333 The odds ratio predicting closely related isolates (defined by having five or fewer single nucleotide  
334 variants between them) associated with sharing a particular lineage (relative to lineage 4), or having  
335 identical or similar MIRU-VNTR profiles, stratified by lineage. PPV denotes positive predictive values.  
336 n refers to the number of subjects having the property described. For example, there were 954  
337 subjects of lineage 4.

338 *Figure 4 The relationship between lineage, MIRU-VNTR profile variation and SNV variation*

339 The relationship between MIRU-VNTR profile variation and SNV variation, stratified by lineage. The  
340 x-axis shown the number of MIRU-VNTR repeats differing between pairs of isolates. For example, if  
341 a sample had a MIRU-VNTR profile of 121, and another 111, locus #2 has reduced in repeat number  
342 by one, which counts as a 1 MIRU-VNTR profile repeat number change. The y-axis shows the  
343 median number of SNV in each of a large number of pairs examined. The blue line reflects the  
344 median pairwise distance within all sampled isolates of each lineage. Red dots are fitted median  
345 values from a multivariable model fitted to MIRU-VNTR profile differences between 0 and 3.

346



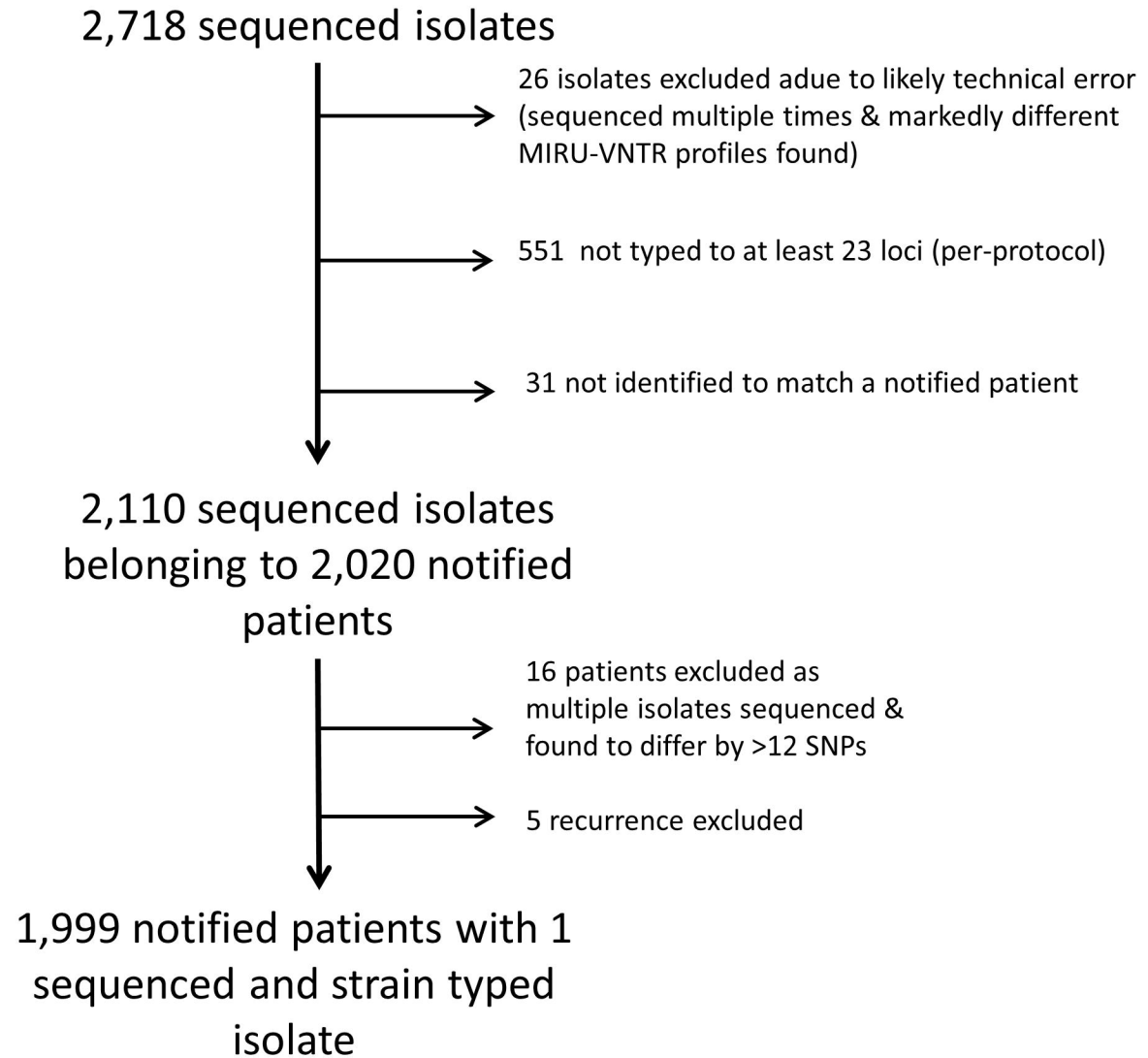
347

348 References

- 349 1. Tuberculosis in England: 2017 Report. London: Public Health England, 2017.
- 350 2. European Centre for Disease Prevention and Control/WHO Regional Office for Europe.  
351 Tuberculosis surveillance and monitoring in Europe 2017. 2017.  
352 [https://ecdc.europa.eu/sites/portal/files/media/en/publications/Publications/ecdc-tuberculosis-](https://ecdc.europa.eu/sites/portal/files/media/en/publications/Publications/ecdc-tuberculosis-surveillance-monitoring-Europe-2017.pdf)  
353 [surveillance-monitoring-Europe-2017.pdf](https://ecdc.europa.eu/sites/portal/files/media/en/publications/Publications/ecdc-tuberculosis-surveillance-monitoring-Europe-2017.pdf) (accessed 1 January 2018 2018).
- 354 3. Lee RS, Behr MA. The implications of whole-genome sequencing in the control of  
355 tuberculosis. *Theor Adv Infect Dis* 2016; **3**(2): 47-62.
- 356 4. Wlodarska M, Johnston JC, Gardy JL, Tang P. A microbiological revolution meets an ancient  
357 disease: improving the management of tuberculosis with genomics. *Clin Microbiol Rev* 2015; **28**(2):  
358 523-39.
- 359 5. Shamputa IC, Jugheli L, Sadradze N, et al. Mixed infection and clonal representativeness of a  
360 single sputum sample in tuberculosis patients from a penitentiary hospital in Georgia. *Respiratory*  
361 *research* 2006; **7**: 99.
- 362 6. Mears J, Abubakar I, Crisp D, et al. Prospective evaluation of a complex public health  
363 intervention: lessons from an initial and follow-up cross-sectional survey of the tuberculosis strain  
364 typing service in England. *BMC public health* 2014; **14**: 1023.
- 365 7. Mears J, Vynnycky E, Lord J, et al. The prospective evaluation of the TB strain typing service  
366 in England: a mixed methods study. *Thorax* 2016; **71**(8): 734-41.
- 367 8. Quan TP, Bawa Z, Foster D, et al. Evaluation of whole genome sequencing for Mycobacterial  
368 species identification and drug susceptibility testing in a clinical setting: a large-scale prospective  
369 assessment of performance against line-probe assays and phenotyping. *Journal of clinical*  
370 *microbiology* 2017.
- 371 9. Jamieson FB, Teatero S, Guthrie JL, Neemuchwala A, Fittipaldi N, Mehaffy C. Whole-  
372 genome sequencing of the Mycobacterium tuberculosis Manila sublineage results in less clustering  
373 and better resolution than mycobacterial interspersed repetitive-unit-variable-number tandem-repeat  
374 (MIRU-VNTR) typing and spoligotyping. *Journal of clinical microbiology* 2014; **52**(10): 3795-8.
- 375 10. Walker TM, Ip CL, Harrell RH, et al. Whole-genome sequencing to delineate Mycobacterium  
376 tuberculosis outbreaks: a retrospective observational study. *The Lancet Infectious diseases* 2013;  
377 **13**(2): 137-46.
- 378 11. Lee RS, Radomski N, Proulx J-F, et al. Reemergence and Amplification of Tuberculosis in the  
379 Canadian Arctic. *Journal of Infectious Diseases* 2015; **211**(12): 1905-14.
- 380 12. Stucki D, Ballif M, Bodmer T, et al. Tracking a tuberculosis outbreak over 21 years: strain-  
381 specific single-nucleotide polymorphism typing combined with targeted whole-genome sequencing. *J*  
382 *Infect Dis* 2015; **211**(8): 1306-16.
- 383 13. Roetzer A, Diel R, Kohl TA, et al. Whole genome sequencing versus traditional genotyping for  
384 investigation of a Mycobacterium tuberculosis outbreak: a longitudinal molecular epidemiological  
385 study. *PLoS Med* 2013; **10**(2): e1001387.
- 386 14. Stucki D, Ballif M, Egger M, et al. Standard Genotyping Overestimates Transmission of  
387 Mycobacterium tuberculosis among Immigrants in a Low-Incidence Country. *Journal of clinical*  
388 *microbiology* 2016; **54**(7): 1862-70.
- 389 15. Lunter G, Goodson M. Stampy: a statistical algorithm for sensitive and fast mapping of  
390 Illumina sequence reads. *Genome Res* 2011; **21**(6): 936-9.
- 391 16. Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and SAMtools.  
392 *Bioinformatics* 2009; **25**(16): 2078-9.
- 393 17. Coll F, McNerney R, Guerra-Assuncao JA, et al. A robust SNP barcode for typing  
394 Mycobacterium tuberculosis complex strains. *Nature communications* 2014; **5**: 4812.

- 395 18. Mazariegos-Canellas O, Do T, Peto T, et al. BugMat and FindNeighbour: command line and  
396 server applications for investigating bacterial relatedness. *BMC bioinformatics* 2017; **18**(1): 477.
- 397 19. Nikolayevskyy V, Kranzer K, Niemann S, Drobniewski F. Whole genome sequencing of  
398 *Mycobacterium tuberculosis* for detection of recent transmission and tracing outbreaks: A systematic  
399 review. *Tuberculosis*; **98**: 77-85.
- 400 20. Comas I, Homolka S, Niemann S, Gagneux S. Genotyping of genetically monomorphic  
401 bacteria: DNA sequencing in *Mycobacterium tuberculosis* highlights the limitations of current  
402 methodologies. *PLoS One* 2009; **4**(11): e7815.
- 403 21. Mears J, Abubakar I, Cohen T, McHugh TD, Sonnenberg P. Effect of study design and setting  
404 on tuberculosis clustering estimates using *Mycobacterium* Interspersed Repetitive Units-Variable  
405 Number Tandem Repeats (MIRU-VNTR): a systematic review. *BMJ open* 2015; **5**(1): e005636.
- 406 22. Bryant JM, Schurch AC, van Deutekom H, et al. Inferring patient to patient transmission of  
407 *Mycobacterium tuberculosis* from whole genome sequencing data. *BMC Infect Dis* 2013; **13**: 110.
- 408 23. Gardy JL, Johnston JC, Ho Sui SJ, et al. Whole-genome sequencing and social-network  
409 analysis of a tuberculosis outbreak. *N Engl J Med* 2011; **364**(8): 730-9.
- 410 24. Casali N, Nikolayevskyy V, Balabanova Y, et al. Evolution and transmission of drug-resistant  
411 tuberculosis in a Russian population. *Nat Genet* 2014; **46**(3): 279-86.
- 412

Fig. 1



# Fig. 2

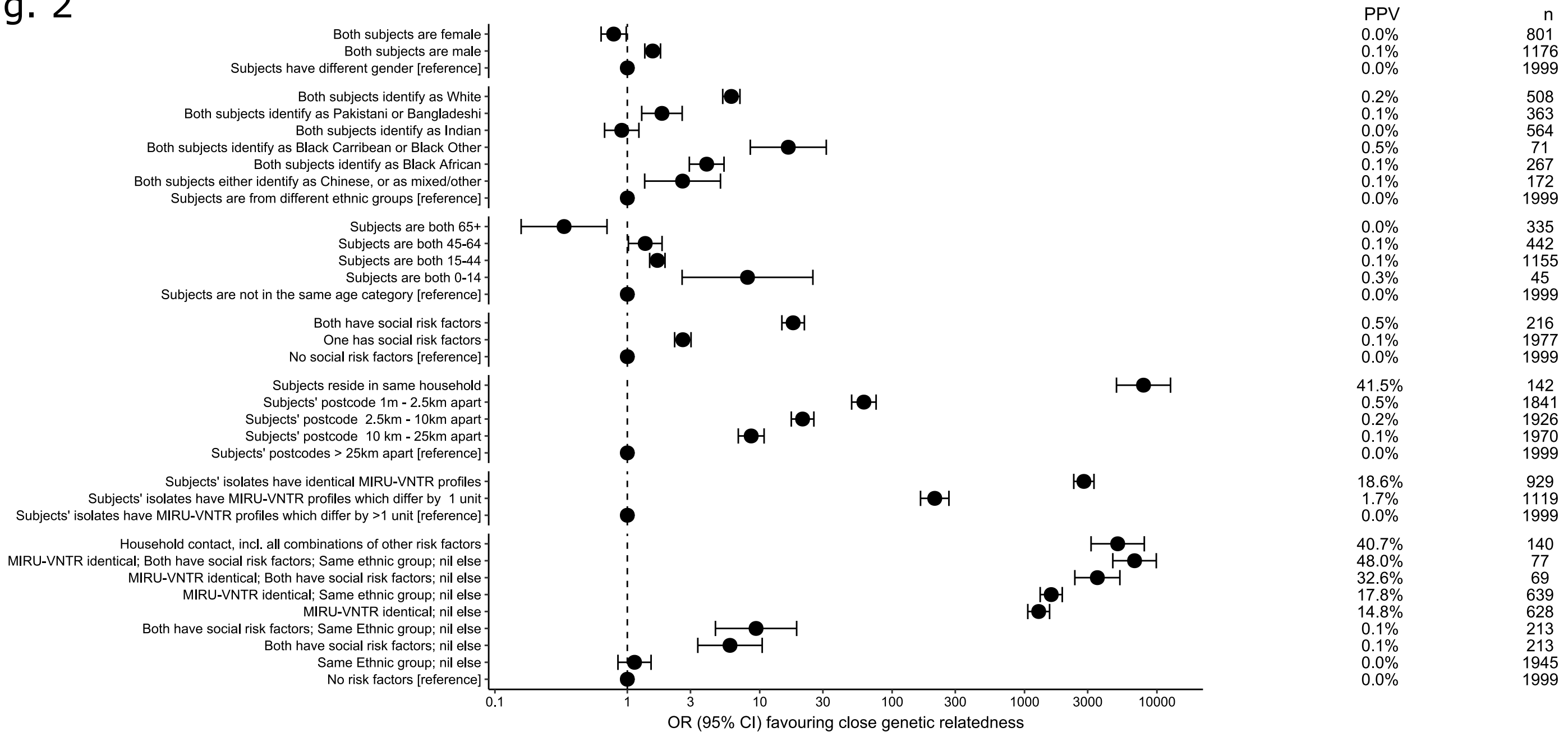


Fig. 3

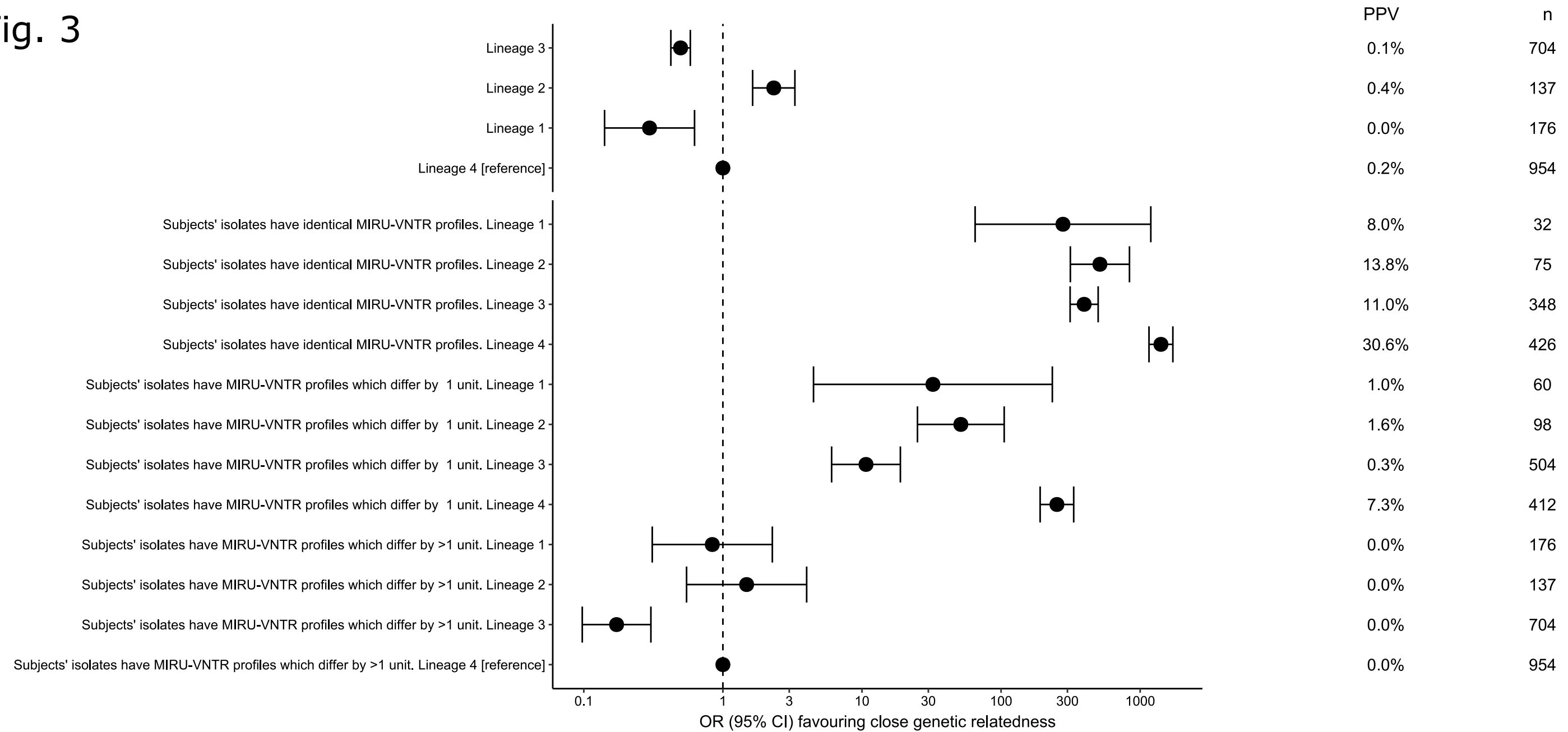


Fig. 4

