

1           Identifying mixed *Mycobacterium tuberculosis* infection and  
2           laboratory cross-contamination during Mycobacterial sequencing  
3           programs

4  
5 Authors

6 David H Wyllie<sup>1, 2, 3 \*</sup>, Esther Robinson<sup>4</sup>, Tim Peto<sup>1, 3</sup>, Derrick W Crook<sup>1, 3</sup>, Adebisi Ajileye<sup>4</sup>, Priti  
7 Rathod<sup>4</sup>, Rosemarie Allen<sup>4</sup>, Lisa Jarrett<sup>4</sup>, E Grace Smith<sup>4</sup>, A Sarah Walker<sup>1, 3</sup>

8  
9 Affiliations

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

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

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

15 <sup>4</sup> Public Health England National Mycobacteriology Laboratory North and Central, Heartlands  
16 Hospital, Birmingham B9 5SS

17

18 Keywords

19 *Mycobacterium tuberculosis*, next generation sequencing, quality control, mixed infection, cross-  
20 contamination

21 ABSTRACT

22 **Introduction:** Detecting laboratory cross-contamination and mixed tuberculosis infection are  
23 important goals of clinical Mycobacteriology laboratories.

24 **Objectives:** To develop a method detecting mixtures of different *M. tuberculosis* lineages in  
25 laboratories performing Mycobacterial next generation sequencing (NGS).

26 **Setting:** Public Health England National Mycobacteriology Laboratory Birmingham, which performs  
27 Illumina sequencing on DNA extracted from positive Mycobacterial Growth Indicator tubes.

28 **Methods:** We analysed 4,156 samples yielding *M. tuberculosis* from 663 MiSeq runs, obtained  
29 during development and production use of a diagnostic process using NGS. Counts of the most  
30 common (major) variant, and all other variants (non-major variants) were determined from reads  
31 mapping to positions defining *M. tuberculosis* lineages. Expected variation was estimated during  
32 process development.

33 **Results:** For each sample we determined the non-major variant proportions at 55 sets of lineage  
34 defining positions. The non-major variant proportion in the two most mixed lineage defining sets  
35 (F2 metric) was compared with that in the 47 least mixed lineage defining sets (F47 metric). Three  
36 patterns were observed: (i) not mixed by either metric, (ii) high F47 metric suggesting mixtures of  
37 multiple lineages, and (iii) samples compatible with mixtures of two lineages, detected by  
38 differential F2 metric elevation relative to F47. Pattern (ii) was observed in batches, with similar  
39 patterns in the H37Rv control present in each run, and is likely to reflect cross-contamination.  
40 During production, the proportions of samples in each pattern were 97%, 2.8%, and 0.001%,  
41 respectively.

42 **Conclusion:** The F2 and F47 metrics described could be used for laboratory process control in  
43 laboratories sequencing *M. tuberculosis*.

44 (249 words)

## 45 INTRODUCTION

46 *Mycobacterium tuberculosis* is an organism which has co-evolved with humans during the early  
47 migrations of modern man, diverging from a common *M. tuberculosis* ancestor about 75,000 years  
48 ago (1). Distinct lineages, corresponding to evolution occurring during these early migrations, are  
49 readily identified by next generation sequencing (NGS), with each lineage characterised by ancient  
50 single nucleotide variants (SNV) which define deep branches in the *M. tuberculosis* phylogeny (1, 2).

51 **Multiple *M. tuberculosis* lineages:** Infection by multiple lineages of TB is well described, and has  
52 been detected by observing mixed results on fractional sequencing (e.g. spoligotyping, MIRU-VNTR)  
53 and validated by characterisation of multiple individual picks from solid media (3). Multi-lineage  
54 infection is characterised by isolates differing by many hundreds of SNVs, in which respect it differs  
55 from the increasingly recognised and more common in-host microevolution (4). Reported rates of  
56 mixed infection vary markedly as reviewed (3, 5), with between 10 and 30% reported in areas of  
57 current (5-7) or historical (8, 9) high prevalence. Much lower rates are reported in low incidence  
58 countries (3), although systematic under-detection is likely to occur due to both to the limited  
59 representation of bacteria in single sputum samples of pulmonary disease, and to the decrease in  
60 diversity occurring during differential strain growth in broth culture (10).

61 **Implications of mixed infection:** Mixed infection, assessed by either MIRU-VNTR polymorphisms (7)  
62 or by heterogeneity in drug susceptibility testing (11), is independently associated with reduced drug  
63 treatment response, so there are compelling clinical reasons to try to identify it. There are also  
64 important technical implications of isolating mixed TB strains from a culture. Firstly, such a finding  
65 may reflect cross-contamination within the laboratory (3). Secondly, mixed infection complicates  
66 the interpretation of drug resistance tests, whether phenotypic or genotypic, as one or other co-  
67 infecting strain may dominate the results from these tests. Thirdly, it complicates the understanding  
68 of relatedness when techniques such as SNV distance computation are applied, as these generally  
69 assume a single sequence is present when base-calling (5, 12-14), marking mixed sites as uncertain.  
70 Maximal likelihood tree drawing algorithms assume such 'uncertain' sites contain no information  
71 and impute a single nucleotide at each of such positions, an approach which may be inappropriate in  
72 the presence of mixtures.

73 Increasingly, NGS based species and resistance determination is becoming routine in  
74 Mycobacteriology laboratories, and has been deployed in reference laboratories in England (15, 16).  
75 As part of the quality control and accreditation of the routine process now operating in these  
76 laboratories, we describe an approach to identifying mixed samples using Illumina next generation

77 sequencing data, illustrating its use by studying over 4,000 consecutive positive cultures from a  
78 single reference laboratory.

79

## 80 METHODS

### 81 Isolation of DNA from Mycobacteria and sequencing

82 This study includes all mycobacteria processed between 01/06/2015 and 30/12/2017 in the Public  
83 Health England Midlands and North reference laboratory, whose catchment is approximately 15  
84 million people, or about one third of England. Clinical specimens were decontaminated and  
85 inoculated into Mycobacterial Growth Indicator Tubes (MGIT) tubes. During the process (Fig. 1)  
86 positive MGIT tubes were batched when they became available, either following growth in the local  
87 laboratory, or following receipt from another laboratory. Positive control samples (H37Rv strain)  
88 were also grown in MGIT cultures. Batches of positive samples were extracted using a manual  
89 process, exactly as described in the Supplementary Methods in (16). Illumina sequencing libraries  
90 were prepared using Nextera XT chemistry from equal amounts of 12 (20.04.15-01.08.15) or 16  
91 (01.08.15-30.12.17) Mycobacterial DNA extracts (16), using manual steps (Table 1). Positive control  
92 DNA (H37Rv, obtained from ATCC) was included as one of the 12 or 16 extracts in all libraries, either  
93 from a contemporaneously extracted broth culture, or from stored DNA. Libraries were loaded into  
94 an Illumina MiSeq instrument and sequenced (16).

### 95 Bioinformatic processing

96 The routine bioinformatics pipeline deployed by Public Health England has been previously  
97 described (15). Briefly, reads were first processed using the Mykrobe predictor tool, which identifies  
98 *Mycobacterium tuberculosis* using species specific k-mers (17). Specimens identified as containing  
99 *M. tuberculosis* were further processed (16), mapped to the H37Rv v2 genome (NC\_000962.2) (18),  
100 as described (16), and vcf files generated using Samtools mPileup, with additional basecalling using  
101 GATK VariantAnnotator v2.1. A consensus base is called from high-quality bases provided one base  
102 accounts from >90% of the pileup; otherwise, the base is recorded as uncertain ('N') (15). In this  
103 analysis, high quality base counts (identified by the BaseCounts VCF tag) were extracted and  
104 summarised using code available at <https://github.com/davidhwylle/VCFMIX>.

### 105 Nucleotides identifying lineage

106 Coll *et al* (2) described the *M. tuberculosis* phylogeny and identified 62 sets of nucleotide positions  
107 defining the deep branches of the *Mycobacterium tuberculosis* lineage. At each position within a  
108 nucleotide set, in one particular clade, one nucleotide is uniquely present (i.e. is not present in any

109 other of the known clades). These sets contain a median 108 nucleotide positions (range 1 – 898).  
110 In this analysis, we considered 55 branches, excluding branches 1.2, 3.1, 3.1.2, 4.1.2, 4.3.4.2.1, 4.6  
111 and 4.7 because they contain fewer than 20 positions, making estimates of minor variation in these  
112 positions less reliable than estimates in other branches. We identified lineage using consensus  
113 basecalling in these 55 branches. If the signature SNV of a branch was called as uncertain, we called  
114 only to the level of the branch deeper (i.e. closer to the root) than the uncertain call. If more than  
115 one different lineage defining variant was called, or we could not call any lineage defining positions,  
116 we reported the samples as ‘lineage not defined’.

#### 117 Estimation of the minor variant frequency within a set

118 The minor variant frequency at a set of bases can be due to sequencing error, mapping error and/or  
119 *bona fide* inter-lineage mixtures (Supplementary Figure 1, panels A and B). Minor variant  
120 frequencies were determined as follows: if there are  $n$  bases in a lineage defining set, we count the  
121 high quality depths  $d$  at each base, e.g. if  $n=3$  and  $d_1 = 30$ ,  $d_2 = 70$ ,  $d_3 = 100$ , then the total depth  
122  $D = \sum_{i=1}^n d_i = 200$ . For each position, we also identify the most common base; the minor depth  $m$   
123 is the total depth minus the most common base depth. If the minor depths are  $m_1 = 3$ ,  $m_2 = 7$ ,  $m_3 =$   
124  $10$ , then total minor depth  $M = \sum_{i=1}^n m_i = 20$ ; we estimate the minor allele fraction  $p$  in the set as  
125  $M/D = 0.1$ .

#### 126 F2 and F47 metrics

127 If sequences from two different *M. tuberculosis* lineages are mixed together, then the sets which  
128 uniquely define these lineages will be mixed (depicted in Supplementary Figure 1C); there will be a  
129 minimum of two and maximum of eight sets affected (e.g. a Lineage 5/7 mixture will mix two sets of  
130 lineage defining nucleotides, a 2.1/4.2.1 will mix five sets, and a 4.1.1.1 / 3.1.2.1 mixture will mix 8  
131 sets). Only if more than two samples are mixed will more than 8 sets be mixed. In this work, we  
132 describe two metrics reflecting mixing. Having computed the minor allele frequency estimates,  $p_1$ ,  
133  $p_2, \dots, p_{55}$  we can sort these in descending order, identifying the sets with the highest and lowest  
134 minor allele frequencies. We then estimate the minor variant frequency across the nucleotides in  
135 the top two (F2 metric) and lowest 47 sets (F47 metric) as described above. The underpinning  
136 assumptions are that mixtures of biological origin are most likely to occur between two lineages, and  
137 therefore F2 is the most sensitive metric for identifying these. Since between two and 8 sets will be  
138 mixed in such genuine co-infections, then the lowest 47 (55 - 8) sets will not be mixed, and thus the  
139 F47 metric is more sensitive for identifying laboratory contamination involving more than two  
140 samples.

141 [Regression modelling](#)

142 Because of high leverage by a small number of observations, we used quantile regression to  
143 estimate the relationship between median values of log-transformed non-callable base numbers and  
144 log-transformed F47, using the quantreg R package (R 3.3.1).

145 [Ethical framework](#)

146 Only anonymised data was used in this work; ethical approval is not required.

147 [Data availability](#)

148 The data analysed is available at [https://ora.ox.ac.uk/objects/uuid:5e4ec1f8-e212-47db-8910-](https://ora.ox.ac.uk/objects/uuid:5e4ec1f8-e212-47db-8910-161a303a0757)  
149 [161a303a0757](https://ora.ox.ac.uk/objects/uuid:5e4ec1f8-e212-47db-8910-161a303a0757).

## 150 RESULTS

### 151 Samples studied

152 4,156 samples were included since they (i) were identified using MyKrobe (17) as belonging to the  
153 *M. tuberculosis* complex, and (ii) had at least  $0.5 \times 10^6$  read pairs mapped to the H37Rv reference  
154 genome, a criterion reflecting successful DNA extraction and sequencing. These sequences were  
155 highly diverse, originating from six branches of lineage 1 (n=320), five branches from lineage 2  
156 (n=278), five branches of lineage 3 (n=1,010), thirty lineage 4 branches (n=2,266). 106 samples were  
157 from *M. bovis* or *africanum*, and 176 did not have their lineages defined.

158 The laboratory processes operated under three different phases: in the first, *development*,  
159 laboratory processes were being actively refined; in the second, *pre-production*, laboratory  
160 processes were fixed and controlled by standard operating procedures, with version controlled  
161 changes. The third *production* state was similar to the second, except that the process had received  
162 ISO15189 accreditation.

### 163 Variation in lineage defining positions in H37Rv controls

164 The F2 mixture metric reflects the estimated mixture in the two most mixed lineage defining sets; in  
165 the H37Rv controls, this follows a distribution skewed to the right (Figure 2A). The MiSeq runs with  
166 H37Rv controls with F2 mixture metrics in the top 5% (Fig. 2A) are temporally clustered (red lines on  
167 Fig. 2B, C) with a number of examples in the Development phase. Among clinical (non-control)  
168 samples, variation in F2 metric is explained in part by MiSeq run (Kruskal-Wallis test,  $p < 10^{-16}$ ), and a  
169 strong correlation exists between the F2 metric in H37Rv controls and that in clinical samples on the  
170 same plate ( $\rho = 0.61$ , 95% CI 0.56-0.61, Spearman's Rank Correlation). That is, in plates with  
171 elevated F2 metrics in the H37Rv control, the clinical samples are more likely to have elevated F2, as  
172 is evident visually (e.g. Fig 2B and C, around run 2301).

### 173 Different patterns of mixtures were observed during Development

174 We ordered specimens first by the order of the plates analysed and the order in which the  
175 bioinformatics processing was completed, which is the order that an automated quality control  
176 monitoring system would encounter output. During the Development phase (Fig. 3), blocks of  
177 samples derived from runs with elevated mixtures in the H37Rv control are seen (red bars in Fig. 3A),  
178 coincident with clear increases in both F2 and F47 metrics (Fig. 3B, C), reflecting elevations in mixed  
179 bases across most or all lineage defining positions (Fig. 3D). These blocks of samples typically span  
180 multiple MiSeq runs (Fig 3A,B). In addition to the blocks of samples with elevated F2 and F47  
181 metrics, we also observed small numbers of single samples with elevated F2, but not F47, metrics  
182 (Fig. 3B, Arrow). This latter pattern is expected in cases of inter-lineage mixtures of only two

183 samples (Supp. Fig 1, and Methods). These patterns were also seen in the subsequent phases  
184 (Supplementary Figures S2-5, yellow dots).

#### 185 [Mixtures of multiple lineages are common](#)

186 Based on the pattern observed in the Development phase, we categorised samples as having one of  
187 (i) neither F2 nor F47 raised, (ii) raised F47, or (iii) raised F2 without raised F47. We defined a raised  
188 F2 and F47 as more than 10x and 5x the respective median metric during Development in all control  
189 and clinical samples, cutoffs which correspond to 4.7% (F2) and 0.2% (F47) minor variant frequencies  
190 across the relevant lineage-defining sets, respectively. In the Pre-production and Production phases,  
191 F2 and F47 values below these thresholds (reflecting unmixed samples) were observed in 97.5% of  
192 the samples studied, raised F47 and F2 values (reflecting a mixture of multiple samples) in 2.5%,  
193 while six samples (0.001%) had raised F2 but normal F47 values (Table 2).

#### 194 [Isolated F2 metric elevation is rare](#)

195 Isolated elevation of F2 is expected if bacteria from two different lineages are mixed. In Fig. 4, we  
196 show the minor variant frequencies from all samples from the six individuals with raised F2, but  
197 normal F47, metrics. In one case, patient 3, two technical repeats of the same sample (#2) showed  
198 the same pattern, as did a separate sample taken contemporaneously. In other cases (patients 1, 5,  
199 6) the mixed pattern was only observed in one out of two positive samples taken on the same day,  
200 and in two cases (2, 3) only a single sample was positive. Thus, between 1 and 6 samples of the  
201 4,156 studied may truly reflect mixed co-infections.

#### 202 [Impact of inter-lineage variation on basecalling](#)

203 One obvious question is whether very low level cross-contamination impacts the consensus  
204 sequence which can be discerned from the pileup. As cross-contamination increases, at some point  
205 minor variant frequencies in some parts of the genome will start to rise above the 10% cutoff  
206 specified by the basecalling algorithm. The numbers of uncalled bases will then rise; this  
207 relationship can be observed in Figure 5: the number of uncalled bases rises rapidly when F47  
208 exceeds the cutoff value, but only slowly below it. Below the cutoff value of 4.7% (red line in Fig. 5),  
209 which is 10 times the median (black line in Fig. 5), the number of uncalled bases increased by 1.25  
210 fold (95% CI 1.22, 1.28) for every 10-fold increase in F47; above the cutoff, the corresponding  
211 increase was 9.24 (95% CI 5.5, 11.2).

212



## 213 DISCUSSION

214 In this work, we describe methods for monitoring the presence of mixtures of *Mycobacterium*  
215 *tuberculosis* of different lineages. It assumes that multiple lineages and sublineages of *M.*  
216 *tuberculosis* are being sequenced contemporaneously; this is the case in our setting, and is also true  
217 globally (19, 20). Using single nucleotide variants, each of which uniquely defines a branch in the  
218 phylogenetic tree of *M. tuberculosis*, we can show two patterns of mixtures. The first, which  
219 occurred in about 2.5% of samples during the pre-production and production phases of our project,  
220 is indicative of multiple samples being mixed together, since mixtures are seen in most or all of the  
221 lineage defining branches. This occurred in batches, was characterised by cross-contamination at  
222 levels of less than 1%, and can be monitored by a metric we term F47. This pattern likely reflects  
223 process failures. The strength of the F47 metric is that the depth analysed is very high, as about  
224 5,000 nucleotides typically contribute across the lineage-defining sets included in it. If there is a  
225 sequencing depth of 50-100 at each of these, the effective sequencing depth analysed is of the order  
226 of 25,000 – 50,000, making detection of minor variation at sub- 1% levels readily feasible with high  
227 statistical confidence.

228 Such low level cross-contamination, as observed during our production process and illustrated in  
229 Supplementary Figures 2-4, is likely to have minimal influence on inference drawn from the  
230 sequence, unless highly sensitive assays for heteroresistance are required. A sensitive metric such as  
231 F47 will, however, allow early detection of emerging problems and allow review of process, as part  
232 of continuous quality improvement.

233 A second class of mixture, which we found was rarely detected in this setting, is compatible with co-  
234 infection with two organisms of differing lineage within the patient. This kind of mixture is clinically  
235 relevant (7, 11), and may be under-detected using the laboratory process we describe here, since  
236 culture based amplification can reduce diversity in the sample inoculated (10). Its frequency may  
237 rise if direct-from-sample sequencing is employed, or if samples from high-endemicity areas are  
238 studied, but here we identified only one probable case of such mixtures, and five other possible  
239 cases. Confirmatory approaches are available: microbiological techniques conducted separately on  
240 multiple picks from the same samples have been used as validation(3); a limitation of this study is  
241 that we could not undertake such work as only multiply sub-cultured stored isolates exist for  
242 historical samples. Techniques for reconstructing the contributing sequences also have been  
243 described in detail (21-23), and we did not study them here. Another limitation is that we were not  
244 able to consider six of the lineage defining sets in Coll *et al* because they covered <20 nucleotide  
245 positions and therefore we considered that they did not contain sufficient information to be used in

246 F2 or F47 metrics. A consequence is that our method would not identify mixtures of samples if they  
247 only involved mixtures in these excluded branches.

248 The clinical of use of bacterial genome sequencing is rising (16, 17, 24), and given the importance of  
249 *M. tuberculosis* and the complexity of treatment, *M. tuberculosis* has been one of the first organisms  
250 tackled (15). The processes followed involve multiple steps at which the opportunity for cross-  
251 contamination exists. The availability of tools monitoring critical aspects of laboratory process is  
252 required for accreditation under ISO15189, and the F2 and F47 metrics described here will  
253 contribute to this.

## 254 [ACKNOWLEDGEMENTS](#)

255 This study is supported by the Health Innovation Challenge Fund (a parallel funding partnership  
256 between the Wellcome Trust [WT098615/Z/12/Z] and the Department of Health [grant HICF-T5-  
257 358]) and NIHR Oxford Biomedical Research Centre. Professors Derrick Crook, Tim Peto and Sarah  
258 Walker are affiliated to the National Institute for Health Research Health Protection Research Unit  
259 (NIHR HPRU) in Healthcare Associated Infections and Antimicrobial Resistance at University of  
260 Oxford in partnership with Public Health England. Professor Crook is based at University of Oxford.  
261 Professor Tim Peto is an NIHR Senior Investigator. The views expressed are those of the author(s)  
262 and not necessarily those of the NHS, the NIHR, the Department of Health or Public Health England.  
263 The sponsors of the study had no role in study design, data collection, data analysis, data  
264 interpretation, or writing of the report. The corresponding author had full access to all the data in  
265 the study and had final responsibility for the decision to submit for publication.

## 266 FIGURE LEGENDS

### 267 Figure 1 Laboratory and Bioinformatic processing

### 268 Figure 2 Mixtures in H37Rv controls

269 Histogram showing the F2 metric, which reflects the mixture in the two most mixed lineage  
270 associated sets, in H37Rv control DNA (A). Median F2 metric among clinical samples other than  
271 H37Rv is shown in (B); red lines indicate that the F2 mixture metric in H37Rv controls is raised (as  
272 shown in A). (C) shows F2 metric for each *M. tuberculosis* sequence from a clinical sample.

### 273 Figure 3 Mixture metrics in the Development phase

274 Samples arranged first by the order of the MiSeq runs (depicted as solid gray blocks, in A), and the  
275 order bioinformatics processing completed. Only samples yielding *M. tuberculosis* are shown, which  
276 is why some blocks in A are longer than others. If the H37Rv control samples had increased F2  
277 statistics, a red bar is shown above each sample in A. We depicted the F2 (B) and F47 metrics (C), as  
278 well as the estimated mixture F in each of the 58 lineage defining sets (D). The arrow illustrates a  
279 sample with elevated F2, but low F47 metric.

### 280 Figure 4 Consistency of isolated F2 elevation in individuals

281 Six individuals with elevated F2, but not F47, statistics were identified during the pre-production and  
282 production phases. The observed minor variant proportion for all deep branches analysed are  
283 shown in a heatmap. For example, patient #4 had two samples taken in December 2015; sample 2  
284 was analysed twice (sequencing ids 0d9d5,9276f) and sample 3 once. A similar pattern of minor  
285 variation is seen in all three samples.

### 286 Figure 5 F47 metric and the number of callable bases

287 The relationship between the F47 metric and the number of uncallable bases. The red line  
288 corresponds to the cutoff used to define F47 as being elevated.

## 289 Supplementary Figures

### 290 *Supplementary Figure 1*

291 Illustrates the computation of variation between samples of two different lineages, Lineage 1 (black  
292 text) and Lineage 2 (blue text) (A). When a mixture of these samples is present, and mapped to a  
293 reference sequence, a major base and minor base(s) are present in the pileup (B). Variation may be  
294 due to either sequencing error (underlined) or to lineage associated variation; the non-major variant  
295 frequency included both classes of variation. Lineage defining sites, as defined by Coll *et al* (2), mark  
296 branches of the phylogenetic tree. If a lineage 4.1.1.1 *M. tuberculosis* is mixed with a lineage 3.1.2.1  
297 *M. tuberculosis*, eight sets of lineage defining sites will be mixed (red boxes).

298 *Supplementary Figures 2-5*

299 These illustrate mixture patterns observed during the Production stage. The layout is similar to  
300 Figure 3; samples arranged first by the order of the MiSeq runs (depicted as solid gray blocks, in A),  
301 and the order bioinformatics processing completed. Only samples yielding *M. tuberculosis* samples  
302 are shown, which is why some blocks in A are longer than others. If the H37Rv control samples had  
303 increase F2 statistics, a red bar is shown above each sample in A. We depicted the F2 metric (B) and  
304 F47 metrics (C), as well as the estimated mixture F in each of the 58 lineage defining sets (D).

305

306

307 TABLES

308 Table 1 Samples analysed

Development Stage	Total sequences (samples and controls)	Date range	Number of individuals providing samples	MiSeq run range	Number of clinical samples	Number of MiSeq Runs
Development	938	20.04.2015-15.12.2015	630	101 .. 291	776	154
Pre-Production	1167	01.04.2016-06.12.2016	753	1152 ..1522	919	163
Production	2191	07.12.16-30.12.2017	1481	1523 .. 2307	1794	346

309

310 Table 2 Detection of mixtures in clinical samples

Development Stage	Neither F2 nor F47 raised	F2 raised, but F47 normal	F47 raised
Pre-Production (n=919)	900 (98%)	5 (0.003%)	14 (1.1%)
Production (n=1794)	1741 (97%)	1 (0.001%)	52 (2.8%)

311

312

313

314

## 315 REFERENCES

- 316 1. Comas I, Coscolla M, Luo T, Borrell S, Holt KE, Kato-Maeda M, Parkhill J, Malla B, Berg S,  
317 Thwaites G, Yeboah-Manu D, Bothamley G, Mei J, Wei L, Bentley S, Harris SR, Niemann S,  
318 Diel R, Aseffa A, Gao Q, Young D, Gagneux S. 2013. Out-of-Africa migration and Neolithic  
319 coexpansion of *Mycobacterium tuberculosis* with modern humans. *Nat Genet* 45:1176-82.
- 320 2. Coll F, McNerney R, Guerra-Assuncao JA, Glynn JR, Perdigao J, Viveiros M, Portugal I, Pain A,  
321 Martin N, Clark TG. 2014. A robust SNP barcode for typing *Mycobacterium tuberculosis*  
322 complex strains. *Nat Commun* 5:4812.
- 323 3. Cohen T, van Helden PD, Wilson D, Colijn C, McLaughlin MM, Abubakar I, Warren RM. 2012.  
324 Mixed-strain *mycobacterium tuberculosis* infections and the implications for tuberculosis  
325 treatment and control. *Clin Microbiol Rev* 25:708-19.
- 326 4. Lieberman TD, Wilson D, Misra R, Xiong LL, Moodley P, Cohen T, Kishony R. 2016. Genomic  
327 diversity in autopsy samples reveals within-host dissemination of HIV-associated  
328 *Mycobacterium tuberculosis*. *Nat Med* 22:1470-1474.
- 329 5. Bryant JM, Harris SR, Parkhill J, Dawson R, Diacon AH, van Helden P, Pym A, Mahayiddin AA,  
330 Chuchottaworn C, Sanne IM, Louw C, Boeree MJ, Hoelscher M, McHugh TD, Bateson AL,  
331 Hunt RD, Mwaigwisya S, Wright L, Gillespie SH, Bentley SD. 2013. Whole-genome sequencing  
332 to establish relapse or re-infection with *Mycobacterium tuberculosis*: a retrospective  
333 observational study. *Lancet Respir Med* 1:786-92.
- 334 6. Wang X, Liu H, Wei J, Wu X, Yu Q, Zhao X, Lyu J, Lou Y, Wan K. 2015. An investigation on the  
335 population structure of mixed infections of *Mycobacterium tuberculosis* in Inner Mongolia,  
336 China. *Tuberculosis (Edinb)* 95:695-700.
- 337 7. Cohen T, Chindelevitch L, Misra R, Kempner ME, Galea J, Moodley P, Wilson D. 2016. Within-  
338 Host Heterogeneity of *Mycobacterium tuberculosis* Infection Is Associated With Poor Early  
339 Treatment Response: A Prospective Cohort Study. *J Infect Dis* 213:1796-9.
- 340 8. Kay GL, Sergeant MJ, Zhou Z, Chan JZ, Millard A, Quick J, Szikossy I, Pap I, Spigelman M,  
341 Loman NJ, Achtman M, Donoghue HD, Pallen MJ. 2015. Eighteenth-century genomes show  
342 that mixed infections were common at time of peak tuberculosis in Europe. *Nat Commun*  
343 6:6717.
- 344 9. Chan JZ, Sergeant MJ, Lee OY, Minnikin DE, Besra GS, Pap I, Spigelman M, Donoghue HD,  
345 Pallen MJ. 2013. Metagenomic analysis of tuberculosis in a mummy. *N Engl J Med* 369:289-  
346 90.
- 347 10. Plazzotta G, Cohen T, Colijn C. 2015. Magnitude and sources of bias in the detection of mixed  
348 strain *M. tuberculosis* infection. *J Theor Biol* 368:67-73.
- 349 11. Zetola NM, Modongo C, Moonan PK, Ncube R, Matlhagela K, Sepako E, Collman RG, Bisson  
350 GP. 2014. Clinical outcomes among persons with pulmonary tuberculosis caused by  
351 *Mycobacterium tuberculosis* isolates with phenotypic heterogeneity in results of drug-  
352 susceptibility tests. *J Infect Dis* 209:1754-63.
- 353 12. Walker TM, Lalor MK, Broda A, Ortega LS, Morgan M, Parker L, Churchill S, Bennett K,  
354 Golubchik T, Giess AP, Del Ojo Elias C, Jeffery KJ, Bowler I, Laurenson IF, Barrett A,  
355 Drobniowski F, McCarthy ND, Anderson LF, Abubakar I, Thomas HL, Monk P, Smith EG,  
356 Walker AS, Crook DW, Peto TEA, Conlon CP. 2014. Assessment of *Mycobacterium*  
357 tuberculosis transmission in Oxfordshire, UK, 2007-12, with whole pathogen genome  
358 sequences: an observational study. *Lancet Respir Med* 2:285-292.
- 359 13. Walker TM, Ip CL, Harrell RH, Evans JT, Kapatai G, Dedicoat MJ, Eyre DW, Wilson DJ, Hawkey  
360 PM, Crook DW, Parkhill J, Harris D, Walker AS, Bowden R, Monk P, Smith EG, Peto TE. 2013.  
361 Whole-genome sequencing to delineate *Mycobacterium tuberculosis* outbreaks: a  
362 retrospective observational study. *Lancet Infect Dis* 13:137-46.
- 363 14. Guerra-Assuncao JA, Houben RM, Crampin AC, Mzembe T, Mallard K, Coll F, Khan P, Banda L,  
364 Chiwaya A, Pereira RP, McNerney R, Harris D, Parkhill J, Clark TG, Glynn JR. 2015. Recurrence  
365 due to relapse or reinfection with *Mycobacterium tuberculosis*: a whole-genome sequencing

- 366 approach in a large, population-based cohort with a high HIV infection prevalence and active  
367 follow-up. *J Infect Dis* 211:1154-63.
- 368 15. Quan TP, Bawa Z, Foster D, Walker T, Del Ojo Elias C, Rathod P, Iqbal Z, Bradley P, Mowbray  
369 J, Walker AS, Crook DW, Wyllie DH, Peto TEA, Smith EG. 2017. Evaluation of whole genome  
370 sequencing for Mycobacterial species identification and drug susceptibility testing in a  
371 clinical setting: a large-scale prospective assessment of performance against line-probe  
372 assays and phenotyping. *J Clin Microbiol* doi:10.1128/jcm.01480-17.
- 373 16. Pankhurst LJ, Del Ojo Elias C, Votintseva AA, Walker TM, Cole K, Davies J, Fermont JM,  
374 Gascoyne-Binzi DM, Kohl TA, Kong C, Lemaitre N, Niemann S, Paul J, Rogers TR, Roycroft E,  
375 Smith EG, Supply P, Tang P, Wilcox MH, Wordsworth S, Wyllie D, Xu L, Crook DW. 2016.  
376 Rapid, comprehensive, and affordable mycobacterial diagnosis with whole-genome  
377 sequencing: a prospective study. *Lancet Respir Med* 4:49-58.
- 378 17. Bradley P, Gordon NC, Walker TM, Dunn L, Heys S, Huang B, Earle S, Pankhurst LJ, Anson L,  
379 de Cesare M, Piazza P, Votintseva AA, Golubchik T, Wilson DJ, Wyllie DH, Diel R, Niemann S,  
380 Feuerriegel S, Kohl TA, Ismail N, Omar SV, Smith EG, Buck D, McVean G, Walker AS, Peto TE,  
381 Crook DW, Iqbal Z. 2015. Rapid antibiotic-resistance predictions from genome sequence  
382 data for *Staphylococcus aureus* and *Mycobacterium tuberculosis*. *Nat Commun* 6:10063.
- 383 18. Lunter G, Goodson M. 2011. Stampy: a statistical algorithm for sensitive and fast mapping of  
384 Illumina sequence reads. *Genome Res* 21:936-9.
- 385 19. Stucki D, Brites D, Jeljeli L, Coscolla M, Liu Q, Trauner A, Fenner L, Rutaihua L, Borrell S, Luo  
386 T, Gao Q, Kato-Maeda M, Ballif M, Egger M, Macedo R, Mardassi H, Moreno M, Tudo  
387 Vilanova G, Fyfe J, Globan M, Thomas J, Jamieson F, Guthrie JL, Asante-Poku A, Yeboah-  
388 Manu D, Wampande E, Ssengooba W, Joloba M, Henry Boom W, Basu I, Bower J, Saraiva M,  
389 Vaconcellos SEG, Suffys P, Koch A, Wilkinson R, Gail-Bekker L, Malla B, Ley SD, Beck HP, de  
390 Jong BC, Toit K, Sanchez-Padilla E, Bonnet M, Gil-Brusola A, Frank M, Penlap Beng VN,  
391 Eisenach K, Alani I, Wangui Ndung'u P, et al. 2016. *Mycobacterium tuberculosis* lineage 4  
392 comprises globally distributed and geographically restricted sublineages. *Nat Genet* 48:1535-  
393 1543.
- 394 20. Manson AL, Cohen KA, Abeel T, Desjardins CA, Armstrong DT, Barry CE, 3rd, Brand J,  
395 Chapman SB, Cho SN, Gabrielian A, Gomez J, Jodals AM, Joloba M, Jureen P, Lee JS, Malinga  
396 L, Maiga M, Nordenberg D, Noroc E, Romancenco E, Salazar A, Ssengooba W, Velayati AA,  
397 Winglee K, Zalutskaya A, Via LE, Cassell GH, Dorman SE, Ellner J, Farnia P, Galagan JE,  
398 Rosenthal A, Crudu V, Homorodean D, Hsueh PR, Narayanan S, Pym AS, Skrahina A,  
399 Swaminathan S, Van der Walt M, Alland D, Bishai WR, Cohen T, Hoffner S, Birren BW, Earl  
400 AM. 2017. Genomic analysis of globally diverse *Mycobacterium tuberculosis* strains provides  
401 insights into the emergence and spread of multidrug resistance. *Nat Genet* 49:395-402.
- 402 21. Eyre DW, Cule ML, Griffiths D, Crook DW, Peto TE, Walker AS, Wilson DJ. 2013. Detection of  
403 mixed infection from bacterial whole genome sequence data allows assessment of its role in  
404 *Clostridium difficile* transmission. *PLoS Comput Biol* 9:e1003059.
- 405 22. Pulido-Tamayo S, Sanchez-Rodriguez A, Swings T, Van den Bergh B, Dubey A, Steenackers H,  
406 Michiels J, Fostier J, Marchal K. 2015. Frequency-based haplotype reconstruction from deep  
407 sequencing data of bacterial populations. *Nucleic Acids Res* 43:e105.
- 408 23. Gan M, Liu Q, Yang C, Gao Q, Luo T. 2016. Deep Whole-Genome Sequencing to Detect Mixed  
409 Infection of *Mycobacterium tuberculosis*. *PLoS One* 11:e0159029.
- 410 24. De Silva D, Peters J, Cole K, Cole MJ, Cresswell F, Dean G, Dave J, Thomas DR, Foster K,  
411 Waldram A, Wilson DJ, Didelot X, Grad YH, Crook DW, Peto TE, Walker AS, Paul J, Eyre DW.  
412 2016. Whole-genome sequencing to determine transmission of *Neisseria gonorrhoeae*: an  
413 observational study. *Lancet Infect Dis* 16:1295-1303.

414

415

# FIGURE 1

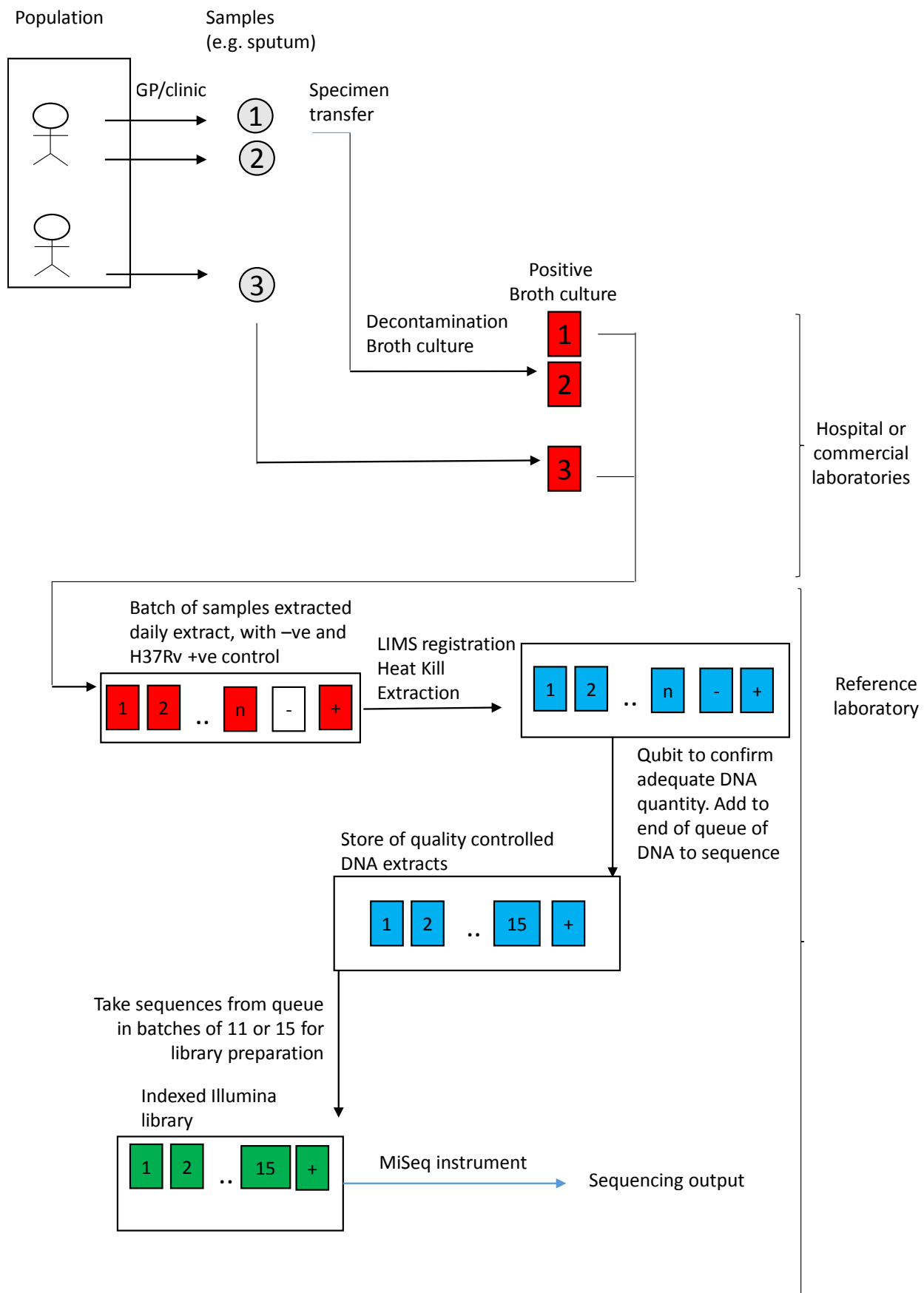




FIGURE 2

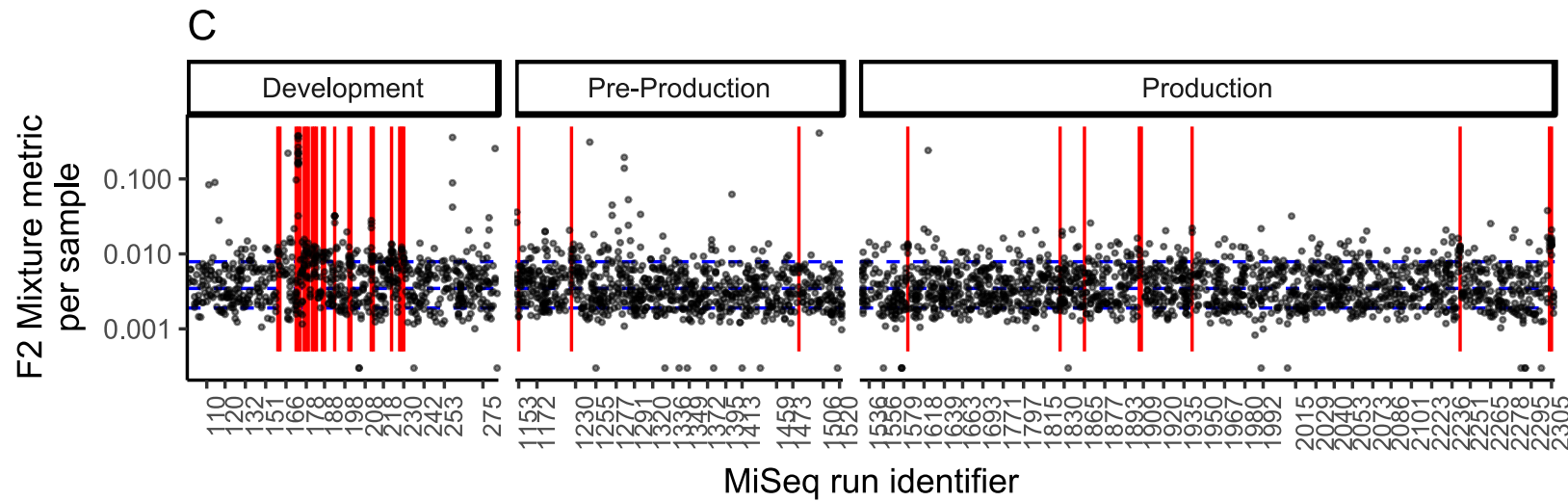
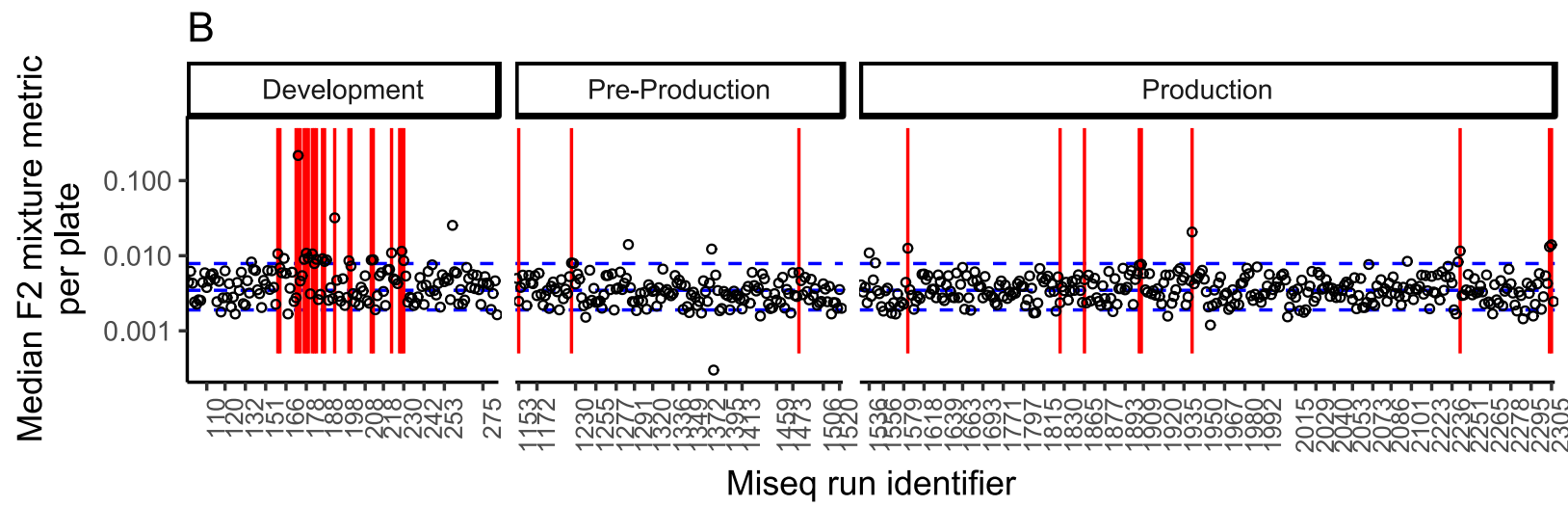
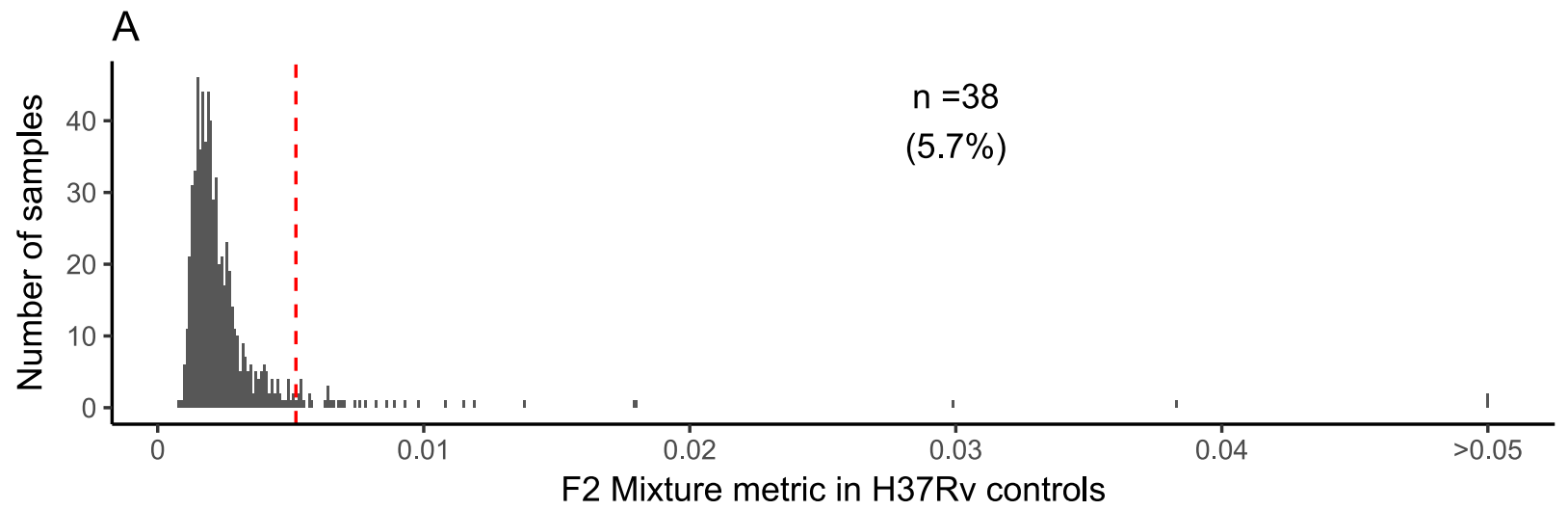


Figure 3

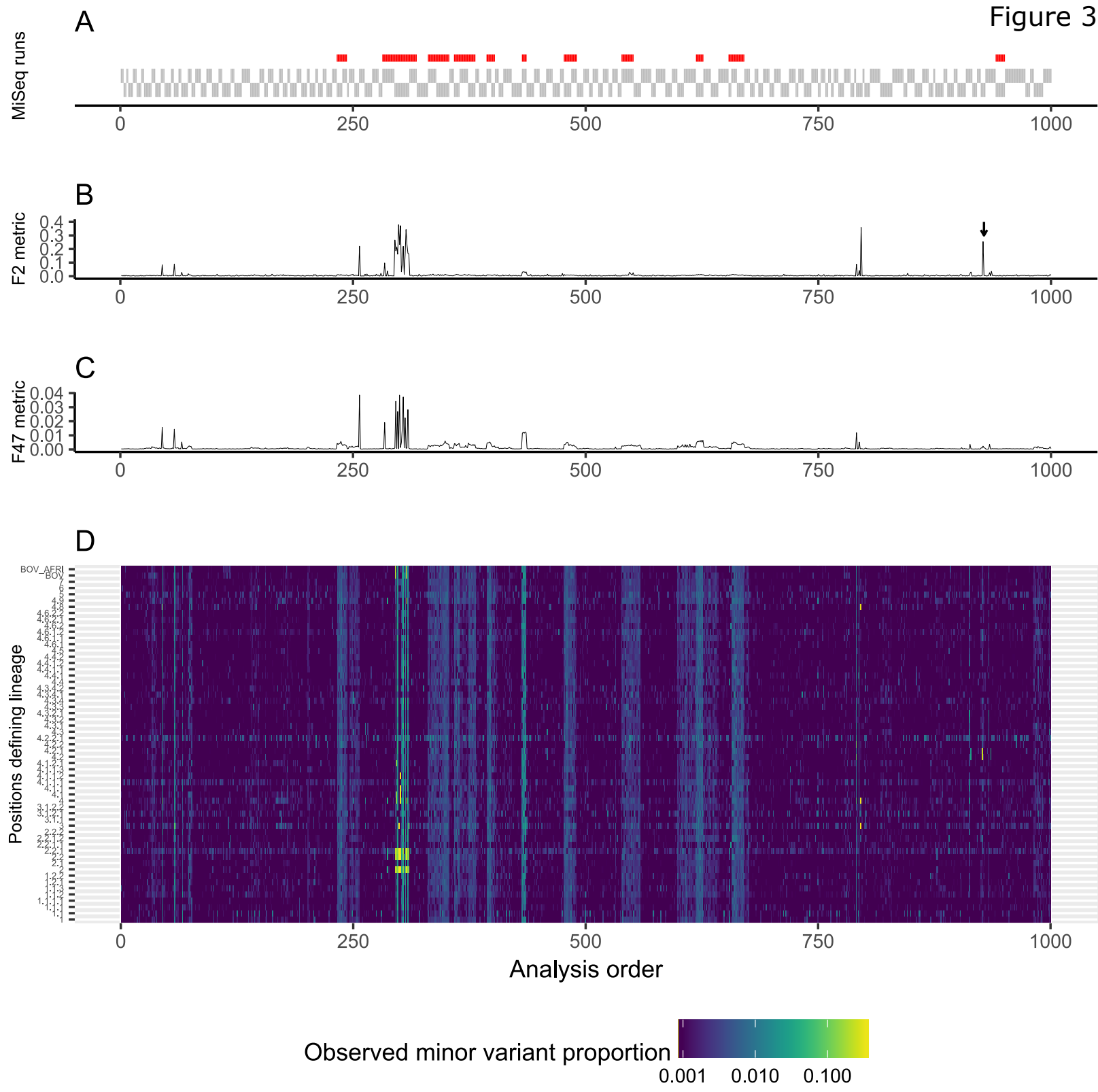
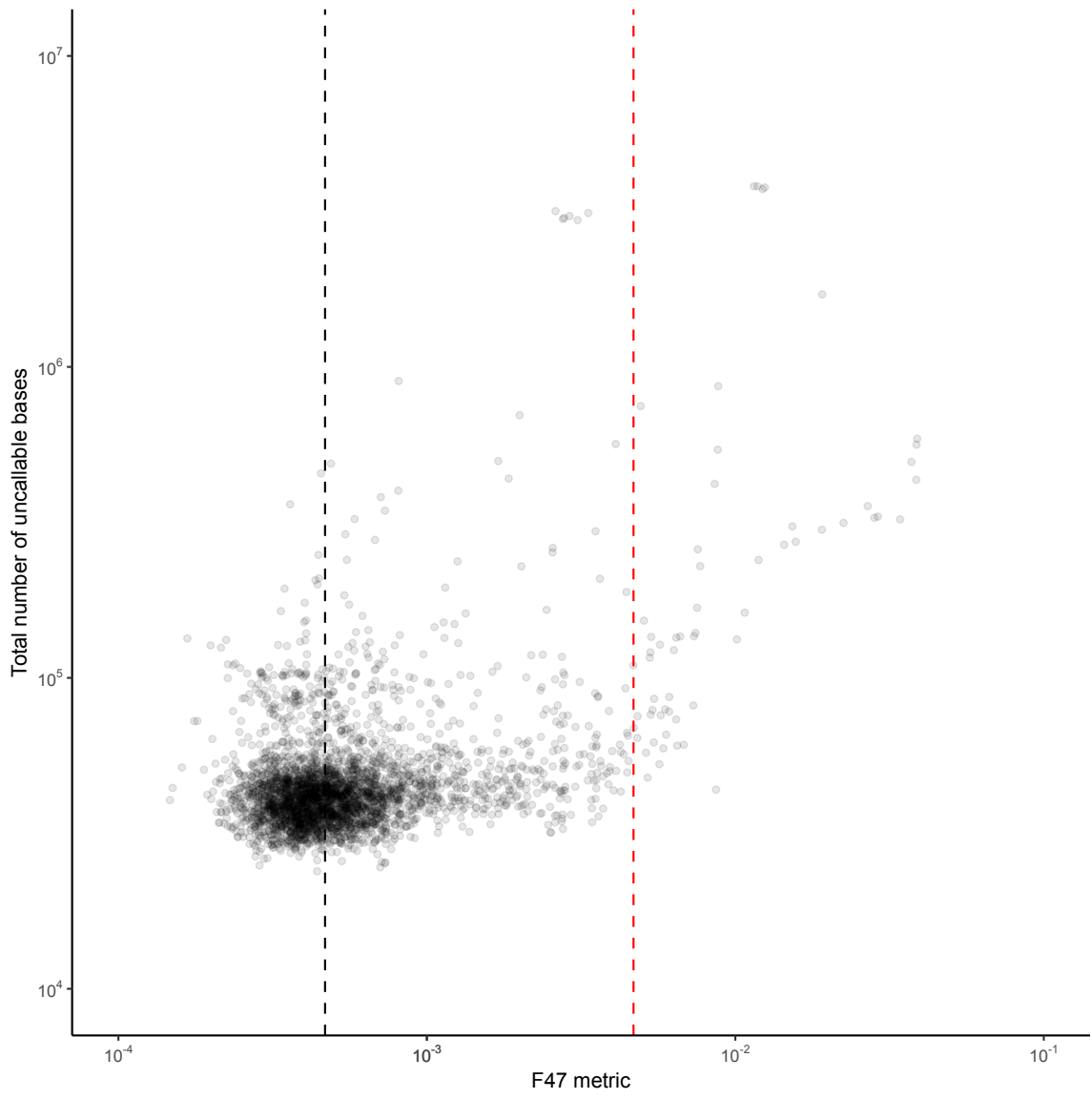




Figure 5



**A**

ACATACGTACGTACGTACGT  
 ACGTACGTT**T**ACGTACGTACGT

Sequence of  
 lineage 1  
 lineage 2

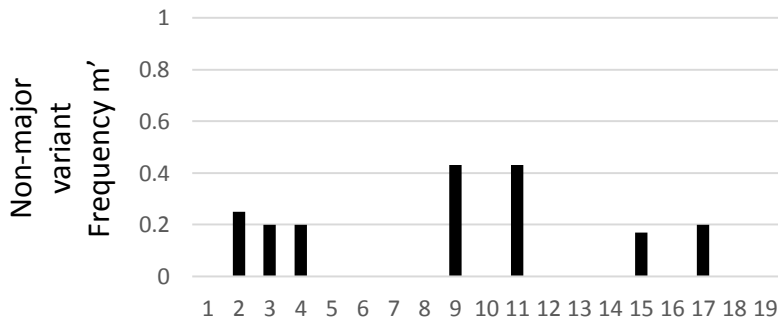
ACATACGTACGT  
 GTACGTACGTACGT  
 CGTACGTACTTACGT  
 GTACGTACGTACGT  
 ACGTACGTT**T**ACGTACGTACGT  
 AAGTACGTT**T**ACGTCCG  
 CGGACGTT**T**ACGTACGTACGT

Reads from  
 lineage 1  
 lineage 2

**T** (bold) lineage  
 defining variant

T (underlined)  
 Variation due to error

**B**



Non-major variant  
 frequency m'

**C**

