

RUNNING TITLE: Deep sequencing as a tool for TB control

1 **Previously undetected superspreading of *Mycobacterium tuberculosis* revealed by**
2 **deep sequencing**

3
4 Robyn S. Lee PhD^{1,2,3*}, Jean-François Proulx MD⁴, Fiona McIntosh BSc⁵, Marcel A. Behr MD⁵,
5 William P. Hanage PhD^{2,3}
6

- 7 1. Epidemiology Division, Dalla Lana School of Public Health, University of Toronto, Health
8 Sciences Building, 155 College Street, Room HS586, Toronto, ON, Canada, M5T 3M7
9 2. Center for Communicable Disease Dynamics, Harvard T. H. Chan School of Public Health,
10 677 Huntington Avenue, Kresge Building, Room 506I, Boston, MA, USA, 02115
11 3. Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA,
12 USA, 677 Huntington Avenue, Kresge Building, 5th floor, Boston, MA, USA, 02115
13 4. Nunavik Regional Board of Health and Social Services, Kuujuaq, Québec, Canada
14 5. The Research Institute of McGill University Health Centre, Montréal, Québec, Canada
15

16 *Address correspondence to:
17 Dr. Robyn Lee, PhD
18 Epidemiology Division
19 Dalla Lana School of Public Health, University of Toronto
20 Health Sciences Building
21 155 College Street, Room HS 586
22 Toronto, Ontario, Canada
23 M5T 3M7
24 Email: robyn.s.c.lee@gmail.com

25 Word count: 3729

26 Figures: 1

27 Tables: 1

28

29

RUNNING TITLE: Deep sequencing as a tool for TB control

30 **Abstract** – 150/150 words

31 Tuberculosis disproportionately affects the Canadian Inuit. To address this, it is imperative we
32 understand transmission dynamics in this population. We investigate whether ‘deep’ sequencing
33 can provide additional resolution compared to standard sequencing, using a well-characterized
34 outbreak from the Arctic (2011-2012, 50 cases). Samples were sequenced to ~500-1000x and
35 reads were aligned to a novel local reference genome generated with PacBio SMRT sequencing.
36 Consensus and heterogeneous variants were identified and compared across genomes. In contrast
37 with previous genomic analyses using ~50x depth, deep sequencing allowed us to identify a
38 novel super-spreader who likely transmitted to up to 17 other cases during the outbreak (35% of
39 all cases that year). It is increasingly evident that within-host diversity should be incorporated
40 into transmission analyses; deep sequencing can facilitate accurately detection of super-spreaders
41 and corresponding transmission clusters. This has implications not only for TB, but all genomic
42 studies of transmission - regardless of pathogen.

RUNNING TITLE: Deep sequencing as a tool for TB control

43 **Introduction.**

44 Tuberculosis (TB) in Canada is highest among the Inuit, an Indigenous population with a rate
45 over 300x that of the non-Indigenous Canadian-born population in 2016.¹ Canada recently set a
46 goal of TB elimination in the Inuit by 2030,¹ which will not be achieved without halting ongoing
47 transmission. Previous studies have used genomic data either alone or in conjunction with
48 classical epidemiology to investigate TB transmission dynamics in the Canadian North,²⁻⁴ with
49 the aim of identifying clusters to help guide public health interventions. Thus far, such studies
50 have relied on identifying consensus single nucleotide polymorphisms (cSNPs), consistent with
51 prevailing methodology in this field.

52
53 Recent studies suggest that incorporation of within-host diversity into genomic analyses may
54 provide greater resolution of transmission than cSNP-based approaches alone.⁵⁻⁸ This may be
55 particularly important for investigation of outbreaks occurring over short time scales and/or in
56 settings such as the Canadian North, where the genetic diversity of circulating strains is
57 especially low. In both of these circumstances, it is common to find many samples separated by
58 zero cSNPs, hindering accurate source ascertainment. To investigate this hypothesis, we used
59 deep sequencing (i.e., to ~10-fold more than standard, or 500-1000x) to re-evaluate transmission
60 in a densely-sampled outbreak in Nunavik, Québec.

61
62 This outbreak, which has been previously described,^{4,9} comprised 50 microbiologically-
63 confirmed cases of TB who were diagnosed in a single Inuit community between 2011-2012 - a
64 rate of 5,359/100,000 for that year. Genomic epidemiology analyses using sequencing depths of
65 ~50x that are standard in such work, identified multiple clusters of transmission in this outbreak,

RUNNING TITLE: Deep sequencing as a tool for TB control

66 ⁴ however, there was insufficient genetic variation detected to infer precise person-to-person
67 transmission events within these subgroups, given the short time frame and low mutation rate of
68 *M. tuberculosis* (~0.5 SNPs/genome/year for Lineage 4¹⁰). In this study, we illustrate how
69 within-host diversity can be incorporated into transmission analyses and in doing so, find new
70 features of the transmission networks in this community, in particular, identifying a previously
71 unrecognized superspreading event. We highlight a potential role for deep sequencing in public
72 health investigations, with implications for TB control in Canada's North as well as other high-
73 transmission environments.

74

75 **Materials and methods.**

76 **Study subjects.** All 50 samples from the 2011-2012 outbreak⁴ were eligible for inclusion, as
77 well as samples from all cases (n=15) diagnosed in same village in the preceding five years
78 (2007 onwards), 13/15 of which were caused by the same strain of *M. tuberculosis* (the 'Major
79 [Mj]-III' sublineage³). There were two episodes of recurrent TB (i.e., where an individual had
80 microbiologically-confirmed TB once, was cured, but developed TB again during the study
81 period); otherwise, all samples are from unique individuals. All cases had pulmonary TB that
82 was Lineage 4 (Euro-American⁴). Cross-contamination was ruled out as described in⁴.

83

84 **DNA extraction and sequencing.** Laboratory methods are described in detail in the
85 **Supplementary Material.** The Illumina HiSeq 4000 was used for paired-end 100bp sequencing.
86 To obtain the target depth of coverage, pooled libraries were run on four independent lanes.

87

RUNNING TITLE: Deep sequencing as a tool for TB control

88 **Bioinformatics.** Quality control of genomic data is described in detail in the **Supplemental**
89 **Material.** Reads were aligned using the Burrows Wheeler Aligner MEM algorithm (v.0.7.15¹¹)
90 to the H37Rv reference (NC_000962.3 in the National Center for Biotechnology Information
91 [NCBI] RefSeq database) and sorted using Samtools (v.1.5¹²). Analyses were later repeated
92 using a local reference genome (described below). Reads with ambiguous mapping were
93 excluded, as were reads with excessive soft-clipping (i.e., more than 20% of read length) based
94 on our previous work.⁶ Duplicate reads were marked using Picard MarkDuplicates (v.2.9.0,
95 <https://broadinstitute.github.io/picard/>) and reads were locally re-aligned around indels using
96 Genome Analysis ToolKit (GATK, v.3.8¹³). All sites were called using GATK's Unified
97 Genotyper algorithm.

98

99 Variants were filtered for quality using custom Python scripts (v.3.6) with the following
100 thresholds: Phred < 50, Root Mean Squared Mapping Quality (RMS-MQ) \leq 30, depth (DP) < 20,
101 Fisher Strand Bias (FS) \geq 60 and read position strand bias (ReadPos) < -8.⁶ cSNPs were
102 classified as positions where \geq 95% of reads were the alternative allele (ALT), hSNPs were
103 classified as positions where > 5% and < 95% of reads were ALT, and positions with the ALT
104 present in < 5% of reads were classified as 'reference'. We also compared inferences of
105 transmission from this analysis to i) when these thresholds were increased to the minimum
106 values among cSNPs in the initial H37Rv analysis, and ii) when cSNPs were classified using a
107 threshold of \geq 99%, and hSNPs were classified when 1% < ALT < 99%, in order to assess the
108 robustness of inferences to different filtering protocols.

109

RUNNING TITLE: Deep sequencing as a tool for TB control

110 Low-quality variants, variants in proline-proline-glutamic acid (PE) and proline-glutamic-
111 acid/polymorphic-guanine-cytosine-rich sequence (PE_PGRS) genes, transposons, phage and
112 integrase, and positions with missing data, were excluded. All samples were drug-susceptible,
113 except for MT-6429, which was rendered resistant to isoniazid by a frameshift deletion at
114 position 1284 in the catalase-peroxidase gene *katG*. As such, positions associated with drug
115 resistance were not masked in this analysis.

116
117 Concatenated cSNP alignments were generated excluding positions with hSNPs. Pairwise cSNP
118 distances between samples were computed using snp-dists (v.0.6, available at
119 <https://github.com/tseemann/snp-dists>). The frequency of hSNPs at each position in the genome
120 was tabulated and hSNPs were reviewed to identify variants shared between samples.

121
122 ***Phylogenetics and clustering.*** Core cSNP alignments were used to generate maximum
123 likelihood trees using IQ-Tree (v.1.6.8¹⁴). Model selection was based on the lowest Bayesian
124 Information Criterion. Hierarchical Bayesian Analysis of Population Structure¹⁵ was run in R
125 (v.3.5.2) to identify clusters. See the **Supplementary Material** for additional detail.

126
127 ***Single Molecule Real-Time (SMRT) sequencing and assembly.*** To examine the influence of
128 potential alignment errors in identification of hSNPs, we used SMRT sequencing with the
129 PacBio RSII platform to create a local reference genome for the outbreak. Sample MT-0080 was
130 chosen for sequencing because this was previously identified as the probable source for as many
131 as 19 of the 50 cases diagnosed in 2011-2012.⁴ A single colony from the culture was selected for

RUNNING TITLE: Deep sequencing as a tool for TB control

132 SMRT sequencing and Illumina MiSeq (for polishing of the long-read assembly). Further detail
133 is provided in the **Supplementary Material**.

134

135 Long-reads were assembled and corrected using Canu (v.1.7.1¹⁶). Pilon (v.1.23¹⁷) was then used
136 to polish the assembly and was re-run until no further corrections were possible. Quast (v.5.0.2,
137 ¹⁸) was used to evaluate assembly quality. RASTtk (v.2.0¹⁹) was used for annotation, to identify
138 regions for masking as previous.

139

140 ***Epidemiological data.*** Epidemiological and clinical data were collected on all cases and contacts
141 using standardized questionnaires, as part of the routine public health response.

142

143 ***Statistical analyses.*** A two-sample test of proportions was used to compare overall proportions
144 across references, and the Wilcoxon Signed Rank test was used to compare paired SNP
145 distances. Analyses were done in Stata (v.15, StataCorp, College Station, TX, USA).

146

147 ***Data availability.*** Sequencing data and the assembly for MT-0080 are available on the NCBI's
148 Sequence Read Archive under BioProject PRJNA549270.

149

150 ***Ethics.*** Ethics approval was obtained from the Institutional Review Board (IRB) of the Harvard
151 T.H. Chan School of Public Health (IRB18-0552) and the IRB of McGill University Faculty of
152 Medicine (IRB A02-M08-18A). All data was analyzed in non-nominal fashion. This study was
153 done with approval of and in collaboration with the Nunavik Regional Board of Health and
154 Social Services.

RUNNING TITLE: Deep sequencing as a tool for TB control

155

156 **Results.**

157 62/65 (95.4%) available TB samples from cases diagnosed between 2007-2012 were successfully
158 sequenced and passed quality control. This included 48/49 (98.0%) of the samples with an
159 identical Mycobacterial Interspersed Repetitive Units Variable Number Tandem Repeats
160 (MIRU-VNTR) pattern during the outbreak year. The remaining three samples could not be re-
161 grown. Reads that were non-MTBC were removed (**Table S1**) and there was no obvious
162 association between percent contamination and hSNP frequency. Epidemiological and clinical
163 data on all outbreak cases are described in ⁴.

164

165 Average genome coverage and depth across the H37Rv reference was 98.64% [SD 0.07%] and
166 714.53 [SD 92.68], respectively. Our primary filtering protocol yielded 51,430 cSNPs and 4,897
167 hSNPs across all individual samples (**Table S2**). Excluding positions that were invariant
168 compared to the reference or where any sample was missing and/or was low-quality resulted in a
169 core alignment of 860 cSNP positions and 136 hSNP positions (note, these are not mutually
170 exclusive, as positions with cSNPs in some samples may have hSNPs in others).

171

172 42 positions had hSNPs that were shared across all 62 samples (**Table 1, Supplementary**
173 **Dataset 1A**). Depth of coverage at these positions was, on average, 39% higher than the average
174 depth across the same sample (SD 36.7%, **Supplementary Dataset 1B**). Along with manual
175 review of alignments (**Figure S1**), this suggested that many of these were false positives,
176 potentially due to alignment error (e.g., from underlying structural variation in our samples
177 compared to the H37Rv reference).

RUNNING TITLE: Deep sequencing as a tool for TB control

178
179 To address this, we generated a local reference genome for the outbreak, MT-0080_PB. Quality
180 metrics for the MT-0080_PB assembly are given in **Table S3**. Compared to H37Rv, mean
181 genome coverage and depth were higher with MT-0080_PB (at 99.33% [SD 0.09%] and 717.07
182 [SD 93.01], respectively), fewer positions were missing/low-quality ($p < 0.00005$, **Table 1**), and
183 overall, fewer variable positions were detected ($p < 0.00005$). While core cSNP distances were
184 similar between samples regardless of the reference (**Table 1**), the number of hSNPs was greatly
185 reduced using MT-0080_PB (**Table S2**); while 4,897 hSNPs were identified across all individual
186 samples using H37Rv, only 125 hSNPs were identified using MT-0080_PB. There were also no
187 hSNPs shared across all 62 samples using MT-0080_PB. Together, these findings support our
188 hypothesis that alignment error is responsible for many of the detected variants, and indicate a
189 local reference is important for accurate identification of hSNPs. All further results presented are
190 based on the MT-0080_PB alignment.

191
192 A maximum likelihood tree was generated from 94 core cSNP positions (excluding sites
193 invariant across all samples and the reference) compared to MT-0080 (**Figure 1A**). Consistent
194 with previous work,⁴ hierBAPS identified two main sub-lineages ('Mj-V' and 'Mj-III' per³),
195 with three sub-clusters (Mj-IIIA/B/C).

196
197 *hSNPs identify super-spreaders and more accurately resolve transmission clusters.* The core
198 cSNPs and hSNPs between samples are shown in **Supplemental Dataset 2A**, with the sub-
199 groups identified in the original analysis indicated. Overlaying hSNPs with the cSNP-based
200 analysis revealed a novel super-spreader (MT-504) in Cluster Mj-IIIB, undetected by genomic

RUNNING TITLE: Deep sequencing as a tool for TB control

201 epidemiology analyses relying on lower sequencing depth.⁴ MT-504 had smear-positive cavitory
202 disease and was diagnosed in late 2011; previous analyses had found this case shared a single
203 cSNP with four other cases diagnosed from March – December of 2012 (position 276,685
204 according to H37Rv / 276,544 in the MT-0080_PB alignment, **Supplementary Dataset 2**).
205 Coupled with epidemiological data on contact (shared attendance at local community ‘gathering
206 houses’, social venues specifically identified by public health during the outbreak), this strongly
207 supported transmission from MT-504 to other members of this subgroup. In contrast, the other
208 subgroup of Mj-IIIB with 13 cases did not share this cSNP. This initially refuted transmission, as
209 we would expect 0 SNPs to accrue in recent transmission given the short time period, low
210 mutation rate of TB, and overall low diversity of strains circulating in the village (**Figure 1B**).
211 Instead, we previously postulated that the first smear-positive case in this subgroup (MT-2474,
212 diagnosed in May 2012) led to the majority of transmission (note, the first smear-negative case in
213 this subgroup was diagnosed in March 2012). However, deep sequencing data suggest otherwise;
214 these data show that MT-504 harboured both the reference (133 reads [19.1%]) allele, present in
215 the subgroup of 13, as well as the alternative allele (563 reads [80.9%]) at this position (**Figure**
216 **1C**). As MT-504 was the first contagious case diagnosed in Mj-IIIB, and all 13 cases in this
217 subgroup had attended or resided in a gathering house (with 9/13 [69.2%] reporting attendance at
218 the same houses as MT-504), this strongly suggests that MT-504 is in fact the most probable
219 source for both subgroups.

220
221 ***hSNP analysis adds support for suspected transmission.*** Sample 68995 and MT-5543 were
222 from 2007, and were the only strains from the Mj-VA sub-lineage in this village. Previous
223 analysis indicated Mj-VA strains from other villages were distantly related,⁴ while these two

RUNNING TITLE: Deep sequencing as a tool for TB control

224 samples were separated from one another by zero core cSNPs. This suggests direct transmission
225 between these historical cases, a hypothesis strongly supported by hSNP analysis, as the samples
226 share hSNPs that are not found in any other sample in the dataset. These hSNPs were present
227 even when highly conservative filtering thresholds were used (**Supplemental Dataset 2B**).
228 Importantly, these hSNPs were not detected when using H37Rv as the reference.

229
230 ***Potential utility for discriminating TB recurrence.*** Six individuals had TB recurrence in 2011-
231 2012. Paired samples were available for two of these (Patient 1: samples MT-5195 in 2007 and
232 MT-1838 in 2012; Patient 2: samples MT-5543 in 2007 and MT-1206 in 2012, **Figure 1A**).
233 cSNP-based analyses suggested their second episodes of TB were due to re-infection with a new
234 strain, rather than relapse with the strain causing their original disease. Investigation of within-
235 host diversity strongly supported this conclusion; using deep sequencing, we verified that there
236 was a single, different strain present at both baseline and their second episodes of TB. There was
237 no evidence for mixed infection at either baseline or second episode with these strains, more
238 definitively ruling out relapse in this low diversity setting (**Supplemental Dataset 2A/B/C**).

239
240 ***Impact of altering cSNP and hSNP thresholds.*** To ensure we were not missing lower frequency
241 variants using the prior cSNP/hSNP thresholds, we re-ran our analysis such that hSNPs were
242 classified when $1\% < \text{ALT} < 99\%$. Quality scores for individual cSNPs and hSNPs are given in
243 **Table S4** and the core cSNP/hSNP alignment is shown in **Supplemental Dataset 2C**. While our
244 primary analysis using a threshold of $> 95\%$ for cSNPs identified a single cSNP (A>G) shared
245 across all samples compared to MT-0080_PB, close examination of the MT-0080 deep
246 sequencing data (obtained using DNA from a sweep of the plate) showed that this sample had

RUNNING TITLE: Deep sequencing as a tool for TB control

247 both alleles at this position, with only the minority ‘A’ allele (33 reads/1189 [2.8%]) isolated for
248 SMRT sequencing. Based on this, we recommend sequencing samples both using a clean sweep
249 (with an alternative sequencing platform) and a single colony pick when generating a reference
250 genome for TB, as using the latter alone may introduce error and affect epidemiological
251 inferences. With this exception, no other informative hSNPs were detected using these
252 thresholds.

253

254 **Discussion.**

255 As the TB epidemic continues among the Canadian Inuit, targeted public health interventions are
256 essential to halt ongoing transmission. In order to do so, it is important that transmission events
257 and associated risk factors are accurately identified. Our previous work suggested that hSNP
258 analysis could enhance resolution of TB transmission⁶. To investigate how this approach could
259 be applied for TB control, we used deep sequencing to re-examine a major TB outbreak in the
260 Canadian Arctic.

261

262 Several recent studies, including work by the authors⁶, have shown that *M. tuberculosis* within-
263 host diversity can be transmitted between individuals^{8,20}. Using deep sequencing data allowed us
264 to better identify this diversity in a Nunavik outbreak compared to previous analyses with
265 standard sequencing depth,^{3,4} and facilitated detection of a novel super-spreader who was likely
266 responsible for ~1/3 of the cases from 2011-2012. Super-spreading has been described in a
267 number of pathogens,²¹ including TB.²² Our findings suggest this can play an important role in
268 driving TB outbreaks. We therefore propose that investigation of within-host diversity is
269 necessary to ensure detection of such super-spreaders in this context, and potentially other high-

RUNNING TITLE: Deep sequencing as a tool for TB control

270 transmission environments as well, in order to accurately identify transmission networks and
271 associated risk factors. We typically only know that a person is a super-spreader retrospectively,
272 i.e., once they have already transmitted to many others and once corresponding genomic data is
273 available. However, if we can identify the characteristics associated with this phenomenon at the
274 population level, this could be used prospectively to predict whether cases are likely to be super-
275 spreaders - as they are diagnosed. This could allow resources to be better allocated, for example,
276 investigation of their contacts could be prioritized and/or targeted screening of the social venues
277 they attend could be rapidly initiated, potentially leading to faster detection of secondary cases
278 and initiation of prophylaxis for new infections. In the case of MT-504, nearly all of the
279 secondary cases had attended the same local community gathering houses as the putative source;
280 this also strongly suggests the importance of these venues in facilitating transmission in this
281 setting.

282
283 Several studies have used genomics to investigate TB recurrence,²³⁻²⁵ however, the methods used
284 to assess for mixed infection at either time point have been inconsistent and may not be sufficient
285 to discriminate recurrence in settings with low strain diversity. In this analysis, we provide proof-
286 of-principle that deep sequencing can potentially help rule out relapse. The distinction between
287 relapse and re-infection is important at individual and population levels; high rates of relapse in a
288 community would indicate a problem with treatment or adherence, potentially warranting
289 changes to clinical management, while re-infection would indicate the need for public health
290 interventions such as activate case finding. Also, individuals in Nunavik who have had prior
291 treatment for active TB disease in the past are also not routinely offered prophylaxis on re-
292 exposure, based on historical data suggesting ~80% protection is afforded by prior infection.²⁶

RUNNING TITLE: Deep sequencing as a tool for TB control

293 The degree to which re-infection drives recurrence in Nunavik is currently unknown, but if re-
294 infection is the primary cause, this clinical practice may need to be re-evaluated. Population-
295 level genomic studies are currently underway to evaluate this.

296
297 To use deep sequencing to investigate within-host diversity, it is critical we minimize false
298 positive hSNPs. We have shown that using a local strain as a reference not only reduces error,
299 but improves detection of epidemiologically-informative variants. Genomic differences between
300 outbreak strains and H37Rv have been previously illustrated by ^{27,28}, with O'Toole *et al.* ²⁸
301 warning that clinical TB strains may be needed to fully detect virulence genes in reference-based
302 analyses. We propose these are also warranted for hSNP analysis. Where possible, we
303 recommend using long-read sequencing to generate complete and local reference genomes.

304
305 Overall, our study has a number of strengths. Firstly, we had access to a densely-sampled
306 outbreak, which was previously investigated using 'standard' sequencing depth and for which
307 detailed epidemiological data was available. This allowed us to readily compare methodological
308 approaches, showing how and when deep sequencing might be beneficial for public health. In
309 doing this, we have identified important methodological considerations for hSNP detection, with
310 implications for transmission analyses, but also potentially for resistance prediction as well.²⁹
311 Finally, the use of long-read data has allowed us to completely assemble a novel TB genome
312 from Nunavik. This will serve as a valuable resource for future studies of transmission in
313 Nunavik (given the low strain diversity in the region ³), as well as other Inuit territories.

314

RUNNING TITLE: Deep sequencing as a tool for TB control

315 A potential limitation of this work is that, given the historical nature of the outbreak, deep
316 sequencing was done using DNA extracted from culture. Due to methodological challenges of
317 sequencing directly from sputum,³⁰⁻³² few studies have examined the effect of culture on genome
318 diversity. A recent study by Shockey *et al*³³. using cSNPs suggests that variation may be lost
319 during the culturing process, however authors did not examine the impact on hSNPs or
320 transmission, and several other studies^{31,32,34} previously found congruent results between cSNP
321 analyses from culture versus raw samples. In terms of hSNPs, Votintseva *et al.* found no
322 difference in the number detected between approaches.³¹ While^{32,34} reported detecting fewer
323 hSNPs with sequencing from culture versus from sputum, in³⁴, the median hSNPs was only 4.5
324 versus 5 hSNPs, respectively – a difference that may not be clinically significant, regardless of
325 statistical significance. Given the inconsistency of results and paucity of data, further study is
326 needed to understand how hSNP diversity may be affected by the culturing process, and to assess
327 whether this affects inferences of transmission. We note that it is likely that enhanced detection
328 of the hSNPs present in sputum would improve the resolution over that which we present in this
329 work.

330

331 Another potential limitation is that, while we can compare the epidemiological inferences made
332 between our previous analysis and our deep sequencing analysis, the bioinformatics pipelines
333 themselves are not directly comparable. Methods to accurately identify hSNPs and incorporate
334 them into transmission analyses are currently an area of active research. We illustrated in our
335 recent paper⁶ that additional steps and strict thresholds must be used to minimize false positive
336 hSNPs, and have conducted the current analysis in consideration of this. However, we note that
337 pre-filtering, our 2015 analysis found that MT-504 had five reference alleles at position 276,685

RUNNING TITLE: Deep sequencing as a tool for TB control

338 in the H37Rv alignment (out of 75) and randomly downsampling the current data to simulate
339 ~50x yielded similar results (5/47 reads at position 276,544 aligned to MT-0080_PB). As most
340 genomic studies of TB employ conservative thresholds of 75-90% allele frequency to classify
341 cSNPs, many bioinformatics pipelines would consider this heterogeneity as potentially suspect at
342 standard sequencing depth. This therefore reinforces the need for greater depth and/or analytic
343 approaches (e.g.,³⁵) to ensure accurate discrimination of sequencing/analytic error from true
344 variation.

345

346 In summary, we have found evidence of mixed variants with important epidemiologic
347 implications that would not have been detected with standard methods and common filtering
348 criteria. This illustrates that genomic methods, while powerful, still require careful interpretation
349 and can still harbor considerable ambiguity when comparing very close links in a transmission
350 chain – a finding whose relevance likely extends far beyond TB, given the increasing number of
351 pathogens undergoing genomic investigation. We demonstrate that deep sequencing can aid in
352 transmission analyses, in particular by allowing accurate identification of TB super-spreaders
353 and key associated epidemiological characteristics. In terms of TB control, this work has
354 important implications for the Canadian North as well as other regions with high TB
355 transmission; as next-generation sequencing becomes a mainstay in public health surveillance, it
356 is critical we recognize the limitations of analyses done using routine sequencing data and
357 consider where and when deep sequencing might be warranted. Although costs continue to
358 decline, we recognize deep sequencing of all samples in an outbreak may not be economically
359 viable for every public health unit. As such, we propose that public health units using routine
360 sequencing for tuberculosis consider – at a minimum – targeted deep sequencing of the more

RUNNING TITLE: Deep sequencing as a tool for TB control

361 contagious (e.g., smear-positive) cases, in lieu of all samples, to help ensure accurate
362 identification of super-spreaders and clusters of transmission. This may help TB control
363 programmes better understand the risk factors for such transmission and enable prioritization of
364 public health resources in future outbreaks.

365

366

367 **Acknowledgements.** We would like to thank the council of the village for their ongoing support
368 of this work. We also acknowledge staff from the Centre Locales des Services Communautaires
369 and the Nunavik Regional Board of Health and Social Services for their hard work during the
370 outbreak. We would like to thank Dr. Hafid Soualhine (currently at the National Reference
371 Centre for Mycobacteriology in the National Microbiology Laboratory, Public Health Agency of
372 Canada) for previous confirmation of the frameshift deletion in *katG* of MT-6429. We also thank
373 Dr. Anders Gonçalves da Silva and Dr. Glen Carter of the Microbiological Diagnostic Unit
374 Public Health Laboratory at the University of Melbourne for their helpful input on high-quality
375 SMRT sequencing. Library preparation and sequencing for all samples was done at the Genome
376 Québec / McGill Innovation Centre and high-performance computing was done using the
377 Odyssey cluster from the Faculty of Arts and Science, Harvard University. This work was
378 supported by an R01 grant from the National Institutes of Health, awarded to WPH
379 (R01AI128344). RSL is also supported by a Fellowship from the Canadian Institutes of Health
380 Research (MFE 152448). MAB holds a Canadian Institutes of Health Research Foundation Grant
381 (FDN-148362).

382

383 **Competing interests.** Authors have no competing interests to declare.

RUNNING TITLE: Deep sequencing as a tool for TB control

384

385 **Contributions.** RSL and WPH conceived and designed the study. RSL designed and did the
386 analyses, made the tables and figures, interpreted the data, and wrote the first draft of
387 manuscript. JFP provided epidemiological data. MAB provided the bacterial samples, as well as
388 laboratory reagents, Biosafety Level 3 access and technician labour in-kind to RSL. FM did the
389 culture and DNA extraction for the HiSeq and PacBio SMRT sequencing. WPH reviewed the
390 initial draft. All authors provided feedback and reviewed the final version of the manuscript
391 before submission. RSL had full access to the data and final responsibility for the decision to
392 submit for publication.

393

394

395

RUNNING TITLE: Deep sequencing as a tool for TB control

396 **References.**

- 397 1. Kanatami IT. Inuit Tuberculosis Elimination Framework. 2018.
- 398 2. Tyler AD, Randell E, Baikie M, et al. Application of whole genome sequence analysis to
399 the study of *Mycobacterium tuberculosis* in Nunavut, Canada. *PLoS One* 2017; **12**(10):
400 e0185656.
- 401 3. Lee RS, Radomski N, Proulx JF, et al. Population genomics of *Mycobacterium*
402 *tuberculosis* in the Inuit. *Proc Natl Acad Sci U S A* 2015; **112**(44): 13609-14.
- 403 4. Lee RS, Radomski N, Proulx JF, et al. Reemergence and amplification of tuberculosis in
404 the Canadian arctic. *J Infect Dis* 2015; **211**(12): 1905-14.
- 405 5. Worby CJ, Lipsitch M, Hanage WP. Shared Genomic Variants: Identification of
406 Transmission Routes Using Pathogen Deep-Sequence Data. *Am J Epidemiol* 2017; **186**(10):
407 1209-16.
- 408 6. Martin MA, Lee RS, Cowley LA, Gardy JL, Hanage WP. Within-host *Mycobacterium*
409 *tuberculosis* diversity and its utility for inferences of transmission. *Microb Genom* 2018; **4**(10).
- 410 7. Meehan CJ, Goig GA, Kohl TA, et al. Whole genome sequencing of *Mycobacterium*
411 *tuberculosis*: current standards and open issues. *Nature Reviews Microbiology* 2019.
- 412 8. Seraphin MN, Norman A, Rasmussen EM, et al. Direct transmission of within-host
413 *Mycobacterium tuberculosis* diversity to secondary cases can lead to variable between-host
414 heterogeneity without de novo mutation: A genomic investigation. *EBioMedicine* 2019; **47**: 293-
415 300.
- 416 9. Lee RS, Proulx JF, Menzies D, Behr MA. Progression to tuberculosis disease increases
417 with multiple exposures. *Eur Respir J* 2016; **48**(6): 1682-9.
- 418 10. Menardo F, Duchene S, Brites D, Gagneux S. The molecular clock of *Mycobacterium*
419 *tuberculosis*. *PLoS Pathog* 2019; **15**(9): e1008067.
- 420 11. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
421 *arXiv* 2013: 1303.3997v2.
- 422 12. Li H, Handsaker B, Wysoker A, et al. The Sequence Alignment/Map format and
423 SAMtools. *Bioinformatics* 2009; **25**(16): 2078-9.
- 424 13. McKenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit: a MapReduce
425 framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010; **20**(9): 1297-
426 303.
- 427 14. Nguyen LT, Schmidt HA, von Haeseler A, Minh BQ. IQ-TREE: a fast and effective
428 stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol Biol Evol* 2015; **32**(1):
429 268-74.
- 430 15. Cheng L, Connor TR, Siren J, Aanensen DM, Corander J. Hierarchical and spatially explicit
431 clustering of DNA sequences with BAPS software. *Mol Biol Evol* 2013; **30**(5): 1224-8.
- 432 16. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and
433 accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res*
434 2017; **27**(5): 722-36.
- 435 17. Walker BJ, Abeel T, Shea T, et al. Pilon: an integrated tool for comprehensive microbial
436 variant detection and genome assembly improvement. *PLoS One* 2014; **9**(11): e112963.
- 437 18. Gurevich A, Saveliev V, Vyahhi N, Tesler G. QUAST: quality assessment tool for genome
438 assemblies. *Bioinformatics* 2013; **29**(8): 1072-5.

RUNNING TITLE: Deep sequencing as a tool for TB control

- 439 19. Brettin T, Davis JJ, Disz T, et al. RASTtk: a modular and extensible implementation of the
440 RAST algorithm for building custom annotation pipelines and annotating batches of genomes.
441 *Sci Rep* 2015; **5**: 8365.
- 442 20. Guthrie JL, Strudwick L, Roberts B, et al. Whole genome sequencing for improved
443 understanding of Mycobacterium tuberculosis transmission in a remote circumpolar region.
444 *Epidemiology and Infection* 2019; **147**.
- 445 21. Stein RA. Super-spreaders in infectious diseases. *Int J Infect Dis* 2011; **15**(8): e510-3.
- 446 22. Kline SE, Hedemark, L.L., Davies, S.F. Outbreak of tuberculosis among regular patrons of
447 a neighborhood bar. *New Engl J Med* 1995; **333**(4): 222-7.
- 448 23. Witney AA, Bateson AL, Jindani A, et al. Use of whole-genome sequencing to distinguish
449 relapse from reinfection in a completed tuberculosis clinical trial. *BMC Med* 2017; **15**(1): 71.
- 450 24. Bryant JM, Harris SR, Parkhill J, et al. Whole-genome sequencing to establish relapse or
451 re-infection with Mycobacterium tuberculosis: a retrospective observational study. *The Lancet*
452 *Respiratory Medicine* 2013; **1**(10): 786-92.
- 453 25. Guerra-Assuncao JA, Houben RM, Crampin AC, et al. Recurrence due to relapse or
454 reinfection with Mycobacterium tuberculosis: a whole-genome sequencing approach in a large,
455 population-based cohort with a high HIV infection prevalence and active follow-up. *J Infect Dis*
456 2015; **211**(7): 1154-63.
- 457 26. Menzies D. Issues in the management of contacts of patients with active pulmonary
458 tuberculosis. *Can J Public Health* 1997; **3**: 197-201.
- 459 27. Roetzer A, Diel R, Kohl TA, et al. Whole genome sequencing versus traditional
460 genotyping for investigation of a Mycobacterium tuberculosis outbreak: a longitudinal
461 molecular epidemiological study. *PLoS Med* 2013; **10**(2): e1001387.
- 462 28. O'Toole RF, Gautam SS. Limitations of the Mycobacterium tuberculosis reference
463 genome H37Rv in the detection of virulence-related loci. *Genomics* 2017; **109**(5-6): 471-4.
- 464 29. Liu Q, Via LE, Luo T, et al. Within patient microevolution of Mycobacterium tuberculosis
465 correlates with heterogeneous responses to treatment. *Sci Rep* 2015; **5**: 17507.
- 466 30. Brown AC, Bryant JM, Einer-Jensen K, et al. Rapid Whole-Genome Sequencing of
467 Mycobacterium tuberculosis Isolates Directly from Clinical Samples. *J Clin Microbiol* 2015; **53**(7):
468 2230-7.
- 469 31. Votintseva AA, Bradley, P., Pankhurst, L., del Ojo Elias, C., Loose, M., Nilgiriwala, K.,
470 Chatterjee, A., Smith, E.G., Sanderson, N., Walker, T.M., Morgan, M.R., Wyllie, D.H., Walker,
471 A.S., Peto, T.E.A., Crook, D.W., Iqbal, Z. Same-day diagnostic and surveillance data for
472 tuberculosis via whole genome sequencing of direct respiratory samples. *J Clin Microbiol* 2017.
- 473 32. Doyle RM, Burgess, C., Williams, R., Gorton, R., Booth, H., Brown, J., Bryant, J.M., Chan,
474 J., Creer, D., Holdstock, J., Kunst, H., Lozewicz, S., Platt, G., Yara Romero, E., Speight, G., Tiberi,
475 S., Abubakar, I., Lipman, M., McHugh, T.D., Breuer, J. Direct Whole-Genome Sequencing of
476 Sputum Accurately Identifies Drug-Resistant Mycobacterium tuberculosis Faster than MGIT
477 Culture Sequencing. *J Clin Microbiol* 2018; **56**(8): e00666-18.
- 478 33. Shockey AC, Dabney J, Pepperell CS. Effects of Host, Sample, and in vitro Culture on
479 Genomic Diversity of Pathogenic Mycobacteria. *Frontiers in Genetics* 2019; **10**.
- 480 34. Nimmo C, Shaw LP, Doyle R, et al. Whole genome sequencing Mycobacterium
481 tuberculosis directly from sputum identifies more genetic diversity than sequencing from
482 culture. *BMC Genomics* 2019; **20**(1): 389.

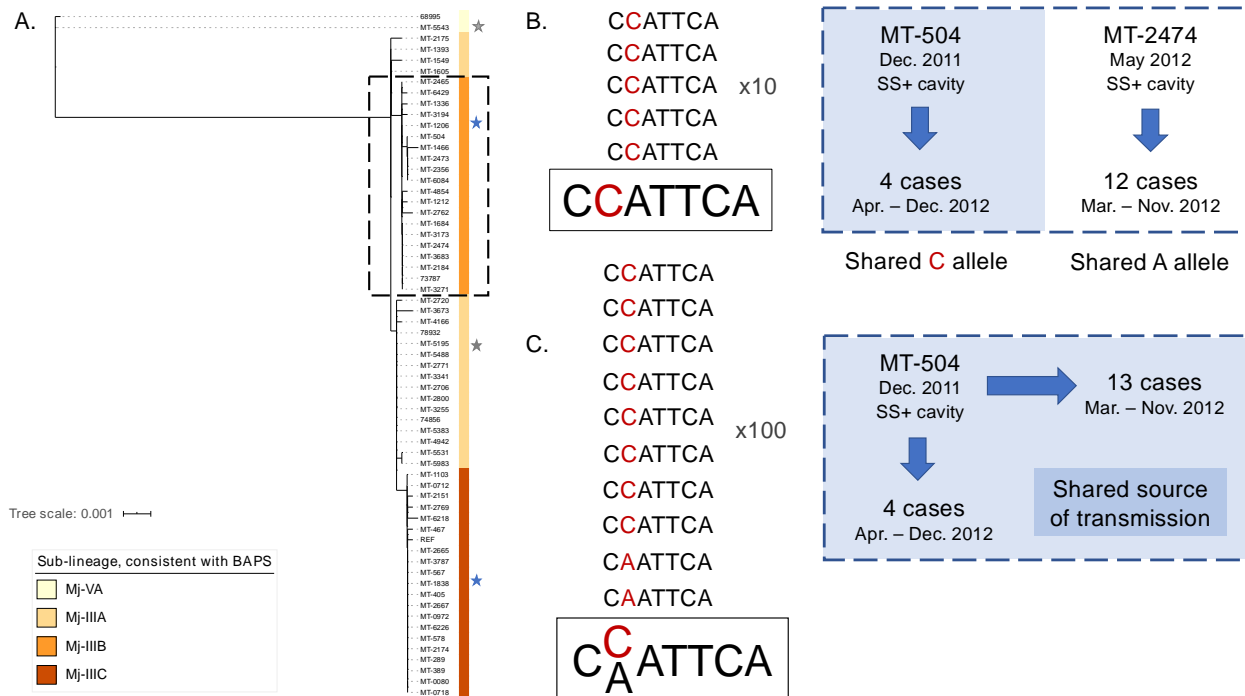
RUNNING TITLE: Deep sequencing as a tool for TB control

483 35. Wyllie D, Do T, Myers R, et al. M. tuberculosis microvariation is common and is
484 associated with transmission: analysis of three years prospective universal sequencing in
485 England. *BioRxiv* 2019.
486

RUNNING TITLE: Deep sequencing as a tool for TB control

487

488 **Figure 1.** Transmission of *M. tuberculosis* in village K.



489

490

491 *Panel A.* Maximum likelihood tree of 62/65 cases diagnosed between 2007-2012 in village K based on consensus
 492 single nucleotide polymorphisms (cSNPs). After aligning to a local reference, MT-0080_PB, cSNPs were identified
 493 based on a minimum threshold of $\geq 95\%$ of reads supporting the alternative allele. A core cSNP alignment was then
 494 produced with 103 positions. IQ-Tree (v.1.6.8¹⁴) was then used to generate the tree using a KP3 model with
 495 correction for ascertainment bias. Model selection was based on the lowest Bayesian Information Criterion. Clusters
 496 were identified using hierarchical Bayesian Analysis of Population Structure.¹⁵ These clusters were consistent with
 497 the sub-lineages previously identified in^{3,4}, thus only sub-lineage names are indicated. During this time period, there
 498 were two individuals who had a second episode of TB; stars are used to highlight these samples, with a different
 499 colour for each patient. MT-0080 is included in the alignment as the deep sequencing data from a sweep of all
 500 colonies identified a cSNP compared to the MT-0080_PB reference, which itself was generated from a single colony
 501 pick.

502

503 *Panel B.* Standard sequencing (to ~40-50x), along with epidemiological data, had indicated that the Major [Mj]-III B
 504 sub-lineage was comprised of two subgroups of five and 13 patients, respectively.⁴ MT-504 was the suspected
 505 source case for the subgroup of five, which all shared a 'C' allele at position 276,685 in H37Rv (position 276,544 in
 506 MT-0080_PB). In contrast, all members of the subgroup of 13 shared an 'A' at this position. Previously, MT-2474
 507 was the suspected source case for this subgroup; this case was the first person with smear-positive (SS+) cavitory
 508 disease diagnosed in this subgroup.

509

510 *Panel C.* In contrast to standard sequencing, deep sequencing data revealed that, in fact, MT-504 – the presumed
 511 source for the subgroup of five cases and the first highly contagious case diagnosed in Mj-III B during the outbreak
 512 year – had both 'C' and 'A' alleles at this position (563:133 of reads, respectively), suggesting this was in fact the
 513 most probable source for both subgroups.

514 **Tables**

515

516 **Table 1.** Comparison of alignments to H37Rv and MT-0080_PB

517

	H37Rv (4,411,532 bp)	MT-0080_PB (4,426,525 bp)	P value
Number of positions according to reference genome			
Invariant reference across all samples, n (%)	4,018,786 (91.10%)	4,084,195 (92.27%)	< 0.00005 ^a
Position was missing / low quality in at least one sample, n (%)	391,761 (8.88%)	342,179 (7.73%)	< 0.00005 ^a
Position was an c/hSNP in at least one sample, n (%)	985 (0.22%)	152 (0.00%)	< 0.00005 ^a
Shared cSNPs across all samples, n (%)	764 (0.02%)	1 (0.00%)	< 0.00005 ^a
Shared hSNPs across all samples, n (%)	42 (0.00%)	0 (0%)	< 0.00005 ^a
Core pairwise distances			
Core cSNPs vs. reference, median (range)	791 (790-792)	3 (1-65)	< 0.00005 ^b
Core cSNPs between samples, median (range)	3 (0-64)	3 (0-66)	< 0.00005 ^b

518 Legend. Based on these filters: Phred < 50, Root Mean Square Mapping Quality (RMS-MQ) ≤ 30, depth (DP) < 20, Fisher Strand Bias (FS) ≥ 60 and read
 519 position strand bias (ReadPos) < -8 and an allelic fraction of 95% for cSNPs, with hSNPs classified when 5% < ALT < 95%. Quality metrics for the individual
 520 cSNPs/hSNPs identified in each sample are given in **Table S2**. ^a Two sample test for difference in proportions. ^b Wilcoxin Signed Rank test.