

1 **Genomic sequence characteristics and the empiric** 2 **accuracy of short-read sequencing**

3
4 Maximillian Marin^{1,2*}, Roger Vargas Jr^{1,2}, Michael Harris³, Brendan Jeffrey³, L. Elaine
5 Epperson⁴, David Durbin⁵, Michael Strong⁴, Max Salfinger⁶, Zamin Iqbal⁷, Irada
6 Akhundova⁸, Sergo Vashakidze^{9,10}, Valeriu Crudu¹¹, Alex Rosenthal³, and Maha Reda
7 Farhat^{1,12*}

8
9 1 Department of Biomedical Informatics, Harvard Medical School, Boston, USA

10 2 Department of Systems Biology, Harvard Medical School, Boston, USA

11 3 Office of Cyber Infrastructure and Computational Biology, National Institute of Allergy and Infectious Diseases,
12 National Institutes of Health, USA

13 4 Center for Genes, Environment, and Health, National Jewish Health, Denver, Colorado, USA

14 5 Mycobacteriology Reference Laboratory, Advanced Diagnostic Laboratories, National Jewish Health, Denver,
15 Colorado, USA

16 6 College of Public Health & Morsani College of Medicine, University of South Florida, USA

17 7 EMBL-EBI, Wellcome Genome Campus, Hinxton, UK

18 8 Scientific Research Institute of Lung Diseases, Ministry of Health, Baku, Azerbaijan

19 9 The University of Georgia, Tbilisi, Georgia

20 10 National Center for Tuberculosis and Lung Diseases, Ministry of Health, Tbilisi, Georgia

21 11 Phthisiopneumology Institute, Chisinau, Moldova

22 12 Pulmonary and Critical Care Medicine, Massachusetts General Hospital, Boston, USA

23 *Corresponding authors: mgmarin@g.harvard.edu, Maha_Farhat@hms.harvard.edu

24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41

42 **Abstract**

43 **Background:** Short-read whole genome sequencing (WGS) is a vital tool for clinical applications
44 and basic research. Genetic divergence from the reference genome, repetitive sequences, and
45 sequencing bias, reduce the performance of variant calling using short-read alignment, but the
46 loss in recall and specificity has not been adequately characterized. For the clonal pathogen
47 *Mycobacterium tuberculosis* (Mtb), researchers frequently exclude 10.7% of the genome believed
48 to be repetitive and prone to erroneous variant calls. To benchmark short-read variant calling, we
49 used 36 diverse clinical Mtb isolates dually sequenced with Illumina short-reads and PacBio long-
50 reads. We systematically study the short-read variant calling accuracy and the influence of
51 sequence uniqueness, reference bias, and GC content. \hat{a}

52 **Results:** Reference based Illumina variant calling had a recall $\geq 89.0\%$ and precision $\geq 98.5\%$ across
53 parameters evaluated. The best balance between precision and recall was achieved by tuning the
54 mapping quality (MQ) threshold, i.e. confidence of the read mapping (recall 85.8%, precision
55 99.1% at MQ ≥ 40). Masking repetitive sequence content is an alternative conservative approach
56 to variant calling that maintains high precision (recall 70.2%, precision 99.6% at MQ ≥ 40). Of the
57 genomic positions typically excluded for Mtb, 68% are accurately called using Illumina WGS
58 including 52 of the 168 PE/PPE genes (34.5%). We present a refined list of low confidence regions
59 and examine the largest sources of variant calling error.

60 **Conclusions:** Our improved approach to variant calling has broad implications for the use of WGS
61 in the study of Mtb biology, inference of transmission in public health surveillance systems, and
62 more generally for WGS applications in other organisms.

63

64 **Background**

65 Illumina short-read whole genome sequencing (WGS) followed by alignment to a reference
66 genome is widely used to identify genetic variants. Illumina sequencing and alignment can
67 confidently detect single nucleotide substitutions (SNSs) and small insertions or deletions (INDELs)
68 but is limited in several ways by its short ~ 100 bp target read lengths. First, short repetitive or
69 homologous query sequences are challenging to uniquely align to the genomic reference^{1,2}.
70 Second, genomic DNA extraction and sequencing library preparation of short-reads may be more
71 error or bias prone³⁻⁷. For example, regions with high GC content and/or low sequence complexity
72 may be particularly prone to PCR-dropout and reduced sequencing coverage⁷⁻⁹. Third, the use of
73 a single reference genome introduces bias, especially when the genome being analyzed differs
74 substantially from the reference sequence^{10,11}. As the sequenced genome diverges from the
75 reference genome, short-read alignment becomes increasingly inaccurate and regions absent
76 from the reference genome are missed or poorly reconstructed.

77

78 In contrast, long-read sequencing can generate high confidence complete genome assemblies,
79 which can also be used to benchmark Illumina WGS. For example, long-reads generated by PacBio
80 sequencing (with lengths on the order of ~10 kb) are ideal for assembling complete bacterial
81 genomes and identifying variants in repetitive regions¹². Although individual PacBio reads have a
82 considerably higher per base error rate (10-15%) than Illumina, the randomly distributed nature
83 of the errors allows for high coverage sequencing runs to converge to a high accuracy consensus¹³.
84 More recently, circular consensus sequencing has further improved PacBio long-read per base
85 accuracy to levels on par with Illumina¹⁴. Alternatively, hybrid strategies that combine less accurate
86 long-reads and short Illumina reads can offer both high base-level accuracy and continuity of the
87 final assembly^{12,15}.

88
89 *Mycobacterium tuberculosis* (Mtb) is a globally prevalent pathogenic bacterium with a ~4.4 Mbp
90 genome known for high GC content, large repetitive regions, and an overall low mutation rate.
91 Owing to the clonality and stability of the Mtb genome, this organism is particularly well suited
92 for systematically identifying the sources of error that arise when short-read data is used for
93 variant detection. Approximately 10% of the Mtb reference genome (H37Rv) is regularly excluded
94 from genomic analysis because it is purported to be more error prone and enriched for repetitive
95 sequence content¹⁶. This 10% of the Mtb genome, hitherto regions of putative low confidence
96 (PLC), span the following genes/families: 1) PE/PPE genes (N=168), 2) mobile genetic elements
97 (MGEs) (N=147), and 3) 69 additional genes with identified homology elsewhere in the genome¹⁷.
98 Despite their systematic exclusion from most Mtb genomic analyses¹⁷⁻¹⁹, PLC regions are yet to
99 be evaluated systematically for short-read variant calling accuracy. Here, we use long-read
100 sequencing data from 36 phylogenetically diverse Mtb isolates to benchmark short-read variant
101 detection accuracy and study genome characteristics that associate with erroneous variant calls.

102

103 **Results**

104 **High confidence Mtb assemblies with hybrid short- and long-read sequencing**

105 For this study, PacBio long-read and Illumina sequencing was performed for 31 clinical Mtb
106 isolates. The resultant data was combined with publicly available paired PacBio and Illumina
107 genome sequencing of 18 Mtb isolates from two previously published studies^{20,21}. From these
108 datasets, a total of 38 clinical isolates were selected for having a) paired end Illumina WGS with
109 median sequencing depth $\geq 40X$ relative to the Mtb reference genome, and b) no evidence of
110 mixed infections or sample swaps (**Additional File 2**).

111

112 Across these 38 isolates, the mean sequencing depth relative to the H37Rv reference genome was
113 84x (IQR: 67x - 107x) for Illumina and 286x (IQR: 180x - 367x) for PacBio. We performed *de novo*
114 genome assembly and iteratively polished each assembly with the PacBio and Illumina reads

115 generating a complete circular assembly for 36/38 isolates (**Methods**). For uniformity in assembly
116 completeness, we excluded the 2 non-circular assemblies from downstream analysis.

117
118 We assessed the accuracy of the *de novo* PacBio assemblies by examining the profile of errors
119 corrected during the Illumina polishing step (**Supp. Figure 1, Additional File 3**). Across all 36
120 assemblies, erroneous 1-bp insertions and deletions (INDELs) made up 97.9% of all corrections
121 made by Illumina polishing with Pilon²². The median number of erroneous insertions and deletions
122 per assembly was 5 (IQR: 2 - 88) and 15 (IQR: 4 - 37) respectively. Very few of the errors corrected
123 during Illumina polishing were single nucleotide changes; median of 0 (IQR: 0 - 2) across all
124 polished 36 genome assemblies. Overall, the number of changes made during Illumina polishing
125 of the *de novo* PacBio assembly was negatively correlated to PacBio sequencing depth
126 (Spearman's $R = -0.458$, $p < 4.9e-3$) (**Supp. Figure 1C**).

127
128 The 36 assemblies spanned the Mtb global phylogeny and had a high degree of conservation in
129 genome structure and content relative to the H37Rv reference genome (**Figure 1, Supp. Figure**
130 **2**): Average Nucleotide Identity (ANI) to H37Rv (99.84% to 99.95%), genome size (4.38-4.44 Mb),
131 GC content (65.59 - 65.64%), and predicted gene count (4017 - 4096 ORFs) (**Additional File 2**).

132
133 In accordance with the small variant benchmarking guidelines of Global Alliance for Genomics &
134 Health²³ (GA4GH), we excluded a small subset of regions with ambiguous ground truths on a per
135 isolate basis (**Methods**). These ambiguous regions fell into 2 categories: a) variable copy number
136 relative to the H37Rv reference genome or b) difficult to align regions due to a high level of
137 sequence divergence relative to the reference genome. We excluded these regions from our
138 performance evaluation in this paper due to their difficulty of interpretation (**Additional File 4**).
139 The percentage of the genome identified as ambiguous was consistently lower than 1% (median:
140 0.41%, IQR: 0.28% - 0.49%) across all assemblies. We observed that for the regions that were
141 frequently ambiguously (Ambiguous in > 25% of isolates, **Additional File 5**), 96.8% of bases were
142 from regions which overlapped with recognized PLC regions.

143

144 **Empirical base-level performance of Illumina**

145 To measure the consistency and accuracy of Illumina genotyping across the Mtb genome, we
146 defined the Empirical Base-level Recall metric (EBR) for each position of the H37Rv reference
147 genome (4.4 Mb, **Additional File 6**). EBR was calculated as the proportion of isolates for which
148 Illumina variant calling made a *confident* variant call that agreed with the ground truth, hence a
149 site with a perfect (1.0) EBR score requires Illumina read data to pass the default quality criteria
150 (**Methods**), and then agree with the PacBio defined ground truth for 100% of the isolates
151 (Examples in **Figure 2**). EBR was significantly lower within PLC regions (mean EBR = 0.905, $N =$
152 469,501 bp) than the rest of the genome (mean EBR = 0.998, $N = 3,942,031$ bp, Mann-Whitney

153 U Test, $P < 2.225e-308$) (**Figure 3A, Table S1**). But EBR was not consistently low across PLC
154 regions, with 67% of PLC base positions having $EBR \geq 0.97$. EBR averaged by gene (gene-level
155 EBR) also showed heterogeneity across PLC regions with 62.6%, 61.3% and 82.6% respectively of
156 the MGEs, PE/PPE, and previously classified repetitive genes having gene-level $EBR \geq 0.97$ (**Figure**
157 **3B, Supp. Figure 3, Tables S2-S3, Additional File 7**). All other, non-PLC, functional gene
158 categories had a median `gene_level_EBR`=1, among these only 14 non-PLC genes had a gene-
159 level $EBR < 0.97$.

160

161 **Characteristics of regions with low empirical performance**

162 Across all 36 isolates evaluated, we observed 1,825,385 sites where Illumina failed to confidently
163 agree with the inferred ground truth. These low recall sites were spread across 267,471 unique
164 positions of the H37Rv reference genome with $EBR < 1$. We explored the underlying factors
165 associated with low recall at these positions using the associated filter and quality tags provided
166 by the variant caller, Pilon (**Methods, Table S4**). Across the 1,829,181 low recall sites, the
167 distribution of outcomes included: a) 62.78% low coverage (LowCov), b) 30.74% falsely called as
168 deleted (Del) with or without low coverage or other tags, c) 6.24% were missed deletions tagged
169 as PASS, d) 0.03% (669 sites) were false base calls (reference or alternate) tagged as PASS, e) 0.25%
170 remaining positions were labeled as ambiguous (Amb) due to evidence for two or more alleles at
171 a frequency $\geq 25\%$.

172

173 Among all low recall sites annotated as with a Low Coverage tag: (a) 45.8% were due to insufficient
174 total coverage of aligned reads (sequencing bias or extreme sequence divergence, total Depth $<$
175 5), (b) 27.6% lacked uniquely aligning reads (repetitive sequence content, mapping quality = 0),
176 and (c) 26.6% were due to low confidence paired-end alignments that did not pass Pilon's
177 heuristics (likely structural variation causing improper paired-alignment orientation).

178

179 **Repetitive sequence content**

180 We identified repetitive regions in H37Rv and evaluated their relationship with low EBR using the
181 pileup mappability metric (**Methods**). Pileup mappability scores range from 0 to 1, where 1
182 represents a genomic position where all overlapping sequence K-mers are unique in the genome
183 of interest within a similarity threshold of E mismatches. We calculated pileup mappability
184 conservatively with a K-mer size of 50 base pairs and up to 4 mismatches (P-Map-K50E4,
185 **Additional File 6**). P-Map-K50E4 is lower in PLC regions (mean = 0.856) than non-PLC regions
186 (mean = .997), (Mann-Whitney U Test, $P < 0.001$) (**Figure 3A**). Yet, 69.7% of positions in PLC
187 regions had P-Map-K50E4 scores of 1, indicating uniquely alignable sequence content even with
188 sequence lengths as short as 50 bp (**Table S5**). At the gene-level, PE/PPEs and MGEs had lower P-
189 Map-K50E4 than the rest of the genome (Wilcoxon, $P < 2e-308$) (**Figure 3B, Table S6, Additional**
190 **File 7**) but 34.5%, and 32.7% of these genes respectively had perfect (1.0) P-Map-K50E4 across

191 the entire gene body. Previously identified repetitive genes (N = 69) had a gene-level P-Map-K50
192 below 1 which is expected given that this was their defining feature²⁴, but for the majority (51 of
193 69), median mappability was greater than 0.99, indicating that a high proportion of their sequence
194 content was actually unique. Non-PLC functional categories had a median gene level P-Map-
195 K50E4 = 1.0 (**Supp. Figure 3, Table S7**). Genome-wide P-Map-K50E4 and EBR scores were
196 moderately correlated (Spearman's $\rho = 0.47$, $P < 2e-308$). Thirty percent of all genome positions
197 with EBR < 1.0 also had a P-Map-K50E4 score below 1.0.

198

199 **Sequencing bias in high GC-content regions**

200 Across several sequencing platforms, high-GC content associates with low sequencing depth due
201 to low sequence complexity, PCR biases in the library preparation and sequencing chemistry³⁻⁶.
202 We assessed the sequencing bias of Illumina and PacBio across each individual genome assembly
203 using the relative depth metric⁴ (the depth per site divided by average depth across the entire
204 assembly) to control for varying depth between isolates. On average with Illumina, 1.2% of the
205 genome had low relative depth (< 0.25), while for PacBio sequencing the average proportion of
206 the genome with low relative depth was 0.0058% (Mann-Whitney U Test, $P < 0.001$). Both
207 sequencing technologies demonstrated coverage bias against high-GC regions, with more
208 extreme bias for Illumina than PacBio (**Figure 4, Additional File 8**). Across all base pair positions
209 with local GC% $\geq 80\%$, using a window size of 100 bp, the mean relative depth was 0.79 for PacBio
210 and 0.35 for Illumina. Genome-wide, EBR was significantly negatively correlated with GC content
211 (Spearman's $\rho = -0.12$, $P < 2e-308$), but this correlation was weaker than that observed with
212 sequence uniqueness (P-Map-K50E4, as above Spearman's $\rho = 0.47$).

213

214 **False positive SNS variant calls**

215 Next, we focused specifically on regions with high numbers of false positive SNSs identified
216 through comparison with the ground-truth variant calls. We examined the distribution of false
217 positive SNS calls across the H37Rv reference genome using a realistic intermediate variant
218 filtering threshold of mean mapping quality at the variant site ($MQ \geq 30$, **Figure 5, Additional**
219 **File 9**). The top 30 regions ranked by the number of false positives (23 genes and 7 intergenic
220 regions) contained 89.4% (490/548) of the total false positive calls and spanned 65 kb, 1.5% of the
221 H37Rv genome. Of these 30 false positive hotspot regions, 29 were either a PLC gene or an
222 intergenic region adjacent to a PLC gene: 17 PE/PPE genes, 3 MGEs, 2 were previously identified
223 repetitive genes²⁴, and 7 PLC-adjacent intergenic regions. Across all false positives, the PE-PGRS
224 and PPE-MPTR sub-families of the PE/PPE genes were responsible for a large proportion (45.4%)
225 of total false positive variant calls. Of all the 556 false positive SNSs evaluated ($MQ \geq 30$), only
226 14 were detected across 4 non-PLC genes: Rv3785 (9 FPs), Rv2823c (1 FP), plsB2 (2 FPs), Rv1435c
227 (2 FPs).

228

229 **Masking to balance precision and recall**

230 A common approach for reducing Mtb false positive variant calls is to mask/exclude all PLC
231 regions from variant calling. Here we investigated two variations on this that utilize directly
232 reference sequence uniqueness and variant quality metrics. We compared: (1) masking of regions
233 with non-unique sequence, defined as positions with P-Map-K50E4 < 1, (2) No *a priori* masking
234 of any regions, and (3) masking of all PLC genes (the current standard practice). We then filtered
235 potential variant calls by whether the variant passed all internal heuristics of the Pilon²²-based
236 variant calling pipeline (**Methods**) and studied the effect of varying the mean mapping quality
237 (MQ) filtering threshold from 1 to 60 (**Figure 6**). We computed the F1-score, precision and recall
238 of detection of SNSs and small indels (<=15bp) for each masking schema and MQ threshold
239 across all 36 clinical isolates (**Methods, Additional File 10**).

240
241 For SNSs, mean recall ranged from 63.6% to 89.0%, and precision ranged from 98.5% to 99.97%
242 across the three schemas (**Figure 6A**). At a threshold of MQ \geq 40, we observed the following mean
243 SNS performances: 1) Masking non-unique regions, F1 = 0.87 (Precision = 99.8%, Recall = 77.9%),
244 2) no masking of the genome, F1 = 0.92 (Precision = 99.1%, Recall = 85.8%), 3) Masking PLC
245 genes, F1 = 0.82 (Precision = 99.6%, Recall = 70.2%). Based on F1 score, no masking of the genome
246 had the highest overall performance, but masking non-unique regions had the highest precision.
247 Decreasing the MQ threshold to an optimal value for F1 score resulted in similar performance for
248 schema-1 and 3, but a balance of lower precision and higher recall for schema-2. Increasing the
249 MQ threshold to 60 optimized precision but at considerable loss of recall for all three schemas
250 (**Table 1**). Performance was most sensitive to the MQ threshold under schema 2 (no masking).

251
252 For INDELs (1-15 bp), precision was comparable to SNSs (96.2% - 100%, **Figure 6B**), while recall
253 was lower (48.9% - 82.4%). At a threshold of MQ \geq 40, we observed the following mean INDEL
254 performances: 1) Masking non-unique regions, F1 = 0.83 (Precision = 98.2, Recall = 72.1%), 2) no
255 masking of the genome, F1 = 0.89 (Precision = 98.9, Recall = 80.8%), 3) Masking PLC genes, F1 =
256 0.76 (Precision = 99.1%, Recall = 61.5%). Variant calling performance of short (1-5bp) INDELs was
257 comparable to SNSs, and the limited performance for INDELs was largely driven by low recall of
258 longer (6-15bp) INDELs (**Supp. Figure 5, Additional File 11**).

259 260 **Structural variation**

261 We assessed the effect of structural variation (SV), of length \geq 50 bp, a common source of
262 reference bias, on variant calling performance (**Methods**). Detected SVs included the known
263 regions of difference associated with Mtb Lineages 1, 2 and 3 (RD239, RD181, RD750
264 respectively)^{25,26} (**Supp. Figure 6**). Across all 36 isolate assemblies, we observed a strong negative
265 correlation between average nucleotide identity to the H37Rv reference and the number of SVs

266 detected (Spearman's $R = -0.899$, $p < 1.1e-13$, **Supp. Figure 7**). Additionally, we observe that 70%
267 of detected SVs overlapped with regions with low pileup mappability ($P\text{-Map-K50E4} < 1.0$).
268

269 We compared SNS variant calling performance by proximity to an SV and sequence uniqueness
270 (**Figure 7, Additional File 12**), dividing variants into four groups: (1) SNSs in regions with perfect
271 mappability ($P\text{map-K50E4} = 1$) with no identified SV (87.3% of total 47,412 SNSs), (2) SNSs in
272 regions with low mappability ($P\text{map-K50E4} < 1$) with no identified SV (10.9% of SNSs), (3) SNSs in
273 regions with perfect mappability within 100 bp of any identified SV (0.8% of SNSs), and (4) SNSs
274 in regions with low mappability within 100bp of any identified SV (1.0% of SNSs). Variant calling
275 performance decreased most sharply in regions with evidence for structural variation, especially
276 when sequence content is also non-unique (Region types 3 & 4 respectively). Additionally, region
277 type (2), or low mappability sequence content with no nearby SV, demonstrated reduced
278 performance.
279

280 **Refined regions of low confidence**

281 Based on the presented analysis, we define a set of refined low confidence (RLC) regions of the
282 Mtb reference genome. The RLC regions are defined to account for the largest sources of error
283 and uncertainty in analysis of Illumina WGS, and is defined as the union of A) The 30 false positive
284 hot spot regions identified (65 kb), B) low recall genomic regions with $\text{EBR} < 0.9$ (142 kb with 30
285 kb overlap with (A)), and C) regions ambiguously defined by long-read sequencing (**Methods**, 16
286 kb). We additionally evaluated the overlap between all detected SVs and the three RLC categories:
287 RLC subset (A) overlapped 28% of SVs, RLC subset (B) overlapped with 65% of SVs, RLC subset (C)
288 overlapped with 14% of SVs.
289

290 In total, the proposed RLC regions account for 177 kb (4.0%) of the total H37Rv genome
291 (**Additional File 13**) and their masking represents a conservative approach to variant filtering.
292 Across the 36 isolates evaluated, masking of the RLC regions combined with a SNS filter of $\text{MQ} \geq$
293 40 would produce a mean F1-score of 0.882, with a mean precision of 99.9% and a mean recall of
294 78.9%.
295

296 **Discussion**

297 The analysis and interpretation of Illumina WGS is critical for both research and clinical
298 applications. Here, we study the 'blindspots' of paired-end Illumina WGS by benchmarking
299 reference-based variant calling accuracy using 36 Mtb isolates with high confidence complete
300 genome assemblies. Overall, our results improve our general understanding of the factors that
301 affect Illumina WGS performance. In particular, we systematically quantify variant calling accuracy
302 and the effect of sequence uniqueness, GC-content, coverage bias, and structural variation. For
303 Mtb, we demonstrate that a much greater proportion of the genome can be analyzed with Illumina

304 WGS than previously thought and provide a systematically defined set of low
305 confidence/troublesome regions for future studies.

306

307 Approaches to benchmarking variant calling from Illumina WGS vary by field and species of
308 interest and more standardization is needed²⁷. Variant calling accuracy is usually benchmarked
309 through *in silico* variant introduction with read simulation or otherwise using a small number of
310 reference genomes that seldom capture the full range of diversity within a particular species. Our
311 benchmarking exercise is unique in using a large and diverse set of high quality genome
312 assemblies that are built using a hybrid long and short read approach. We further demonstrate
313 that PacBio long-read sequencing is much less prone to coverage bias and is able to generate
314 complete circular bacterial assemblies bridging repetitive regions in the majority of isolates with
315 a median depth > 180x. The assemblies we generate will be an important community resource for
316 benchmarking future variant calling or other WGS based bioinformatics tools.

317

318 The benchmarking results clearly demonstrate that low variant recall is a major limitation of
319 reference-based Illumina variant calling, which achieved at most 89% recall at the optimal F1-
320 score. Precision of variant calling using Illumina on the other hand was very high, with the small
321 number of false variant calls concentrated in repetitive and structurally variable regions. We find
322 that the best balance between precision and recall is achieved by tuning the variant mean
323 mapping quality threshold, i.e. confidence of the read mapping. The specific mapping quality
324 threshold will likely vary by species. For a GC-rich organism with highly repetitive sequence
325 content like Mtb, a threshold of 40 achieved 85.8% recall and 99.1% precision.

326

327 Studying specific sources of low recall from Illumina, we identified insufficient read coverage to
328 be the major driver, due not only to repetitive sequence content but also due to high-GC content
329 and other sources of coverage bias. We further identified regions near structural variation to be
330 particularly prone to low recall and precision. Of the variants we study, longer INDELS were recalled
331 at lower rates than SNSs or INDELS < 6bp in length. These observations support ongoing efforts
332 by the bioinformatics research community to build graph-reference genomes and align short
333 reads to these graphs. Using a graph pan-genome built with a diverse set of Mtb reference
334 genomes, there is great potential to both increase recall and precision of variant calling in
335 divergent regions of the genome.

336

337 An alternative and generalizable approach to balancing precision and recall of reference-based
338 Illumina variant calling is to mask repetitive (low mappability) regions. This simple approach does
339 not require tuning the mapping quality threshold against a ground truth set of assemblies and
340 relies instead on computing the pileup mappability metric across the reference sequence. This fills
341 a gap for variant calling in other organisms using short-read mapping where low confidence
342 regions may not already be defined. Compared with tuning against a ground-truth set of

343 assemblies, this masking approach is conservative: for Mtb and filtering by $MQ \geq 40$, precision is
344 slightly higher at 99.8% vs 99.1% respectively and recall is lower at 77.9% vs 85.8% respectively.

345

346 Given Mtb's genomic stability and clonality, this organism is particularly well suited for
347 systematically identifying the sources of variant calling error from short-read data. Although
348 10.7% of the Mtb reference sequence is commonly excluded from genomic analysis, our results
349 demonstrate that more than half of these regions are accurately called using Illumina WGS. For
350 the PE/PPE family, of highest concern for sequencing error, nearly one third (52/168) had perfect
351 mappability and near perfect gene-level EBR (≥ 0.99). The PE/PPE genes with poor performance
352 were largely the PE_PGRS and PPE_MPTR sub-families. Only 65 kb (1.5%) of the reference genome
353 H37Rv were responsible for the majority of false positives (89.2% of false positives across 36
354 isolates).

355

356 We present a set of refined low confidence (RLC) regions of the Mtb genome, designed to account
357 for the largest sources of error and uncertainty in analysis of Illumina WGS (**Additional File 13**).
358 Long-read data can allow RLC regions to be defined for other species to improve accuracy of
359 Illumina WGS. The Mtb RLC regions span 4.0% of the reference genome, and their masking
360 provides a conservative approach to variant calling, appropriate for applications where precision
361 is prioritized over recall. At the same time, RLC region masking offers higher recall than the current
362 field standard where more than 10% of the Mtb reference genome is masked. One limitation is
363 that RLC regions were largely defined based on EBR of Illumina sequencing in our dataset that
364 was restricted by design to 100+ bp paired end sequencing. We do not recommend the use of
365 these RLC regions for Illumina sequencing at shorter read lengths or single-end reads. Instead we
366 make available a more appropriate masking scheme of RLC regions + low pileup mappability
367 (**Additional File 14**). Another limitation is that we defined RLC regions using the same set of high
368 confidence assemblies evaluated. The reported precision and recall with RLC region masking are
369 thus likely overestimates. On the other hand, we expect precision and recall estimates of the
370 alternative approaches of masking low mappability regions or filtering at $MQ \geq 40$ to be more
371 robust.

372

373 Improving Illumina variant recall has significant implications. For clonal Mtb, for example,
374 transmission inference using genomic data often relies on a very small number of SNS or INDEL
375 differences between genome pairs. The observed large increase in recall we observe has the
376 potential to substantially improve transmission inference²⁸ and/or our understanding of genome
377 stability and adaptation.

378

379 **Conclusions**

380 In summary, we show that Illumina whole genome sequencing has high precision but limited recall
381 in repetitive and structurally variable regions when benchmarked against a diverse set of complete
382 assemblies. We demonstrate that filtering variants using the `_mean` mapping quality against a
383 achieves the highest balance of precision and recall. Masking repetitive sequence content is a
384 second generalizable solution, albeit a more conservative one, that maintains high precision. For
385 *Mtb*, these two approaches increase recall of variants by 15.6% and 7.7% respectively, with a
386 minimal change in precision (-0.5% and +0.1% respectively at $MQ \geq 40$), allowing high variant
387 recall in >50% of regions previously considered by the field to be error-prone. Our results improve
388 variant recall from Illumina data with broad implications for clinical and research applications of
389 sequencing. We also provide a high-quality set of genome assemblies for benchmarking future
390 variant calling or other WGS based bioinformatics tools.

391 **Methods**

392 **Summary of sequencing data used**

393 Our dataset consisted of a convenience set of 16 clinical isolates from Lima, Peru, previously
394 sequenced with Illumina WGS and archived in frozen culture²⁹. These isolates were revived and
395 sequenced with PacBio RS II long-read sequencing (Dataset #1). Additionally, 15 total clinical
396 isolates isolated in Azerbaijan, Georgia, Moldova were sequenced with PacBio Sequel II long-read
397 sequencing³⁰ (Dataset #2).

398
399 This dataset of 31 clinical isolates was combined with publicly available paired PacBio (RS II) and
400 Illumina genome sequencing from 19 clinical isolates from two previously published studies^{20,21}.
401 From these four sources, 38 Mtb isolates were selected for having a) Illumina WGS with paired
402 end reads with at least a median sequencing depth of 40X relative to the Mtb reference genome
403 (H37Rv). All aggregated metadata and SRA/ENA accessions for PacBio and Illumina sequencing
404 data associated with this analysis can be found in **Additional File 15**.

405

406 **DNA extraction for PacBio (RS II) Sequencing of Peruvian Isolates (Data Source #1)**

407 MTB cultures were allowed to grow for 4-6 weeks. Pellets were heat-killed at 80°C for 20
408 minutes^{67,68}, the supernatants were removed, and the enriched cell pellet was subjected to DNA
409 extraction soon after or stored frozen until extraction. Largely intact DNA was extracted from heat-
410 killed cells pellets using a protocol tailored for mycobacteria that ends with a column-based
411 elution³¹. Yields were determined using fluorescent quantitation (Qubit, Invitrogen/Thermo Fisher
412 Scientific) and quality was assessed on a 0.8% GelRed agarose gel with 1XTAE, separated for 90
413 minutes at 80V.

414

415 **PacBio (RS II) Sequencing of Peruvian Mtb Isolates (Data Source #1)**

416 Approximately 1 µg of high molecular weight genomic DNA was used as input for SMRTbell
417 preparation, according to the manufacturer's specifications (SMRTbell Template Preparation Kit
418 1.0, Pacific Biosciences). Briefly, HMW gDNA was sheared to 20kb using the Covaris g-tube at 4500
419 rpm. Following shearing, gDNA underwent DNA damage repair, ligation to SMRTbell adaptors
420 and exonuclease treatment to remove any unligated gDNA. At least 500 ng final SMRTbell library
421 per sample was cleaned with AMPure PB beads and 3-50 kb fragments were size selected using
422 the BluePippin system on 0.75% agarose cassettes and S1 ladder, as specified by the manufacturer
423 (Sage Science). Size selected SMRTbell libraries were annealed to sequencing primer and bound
424 to the P6 polymerase prior to loading on the RSII sequencing system (Pacific Biosciences).
425 Sequencing was performed using C4 chemistry and 240-minute movies. Following data collection,
426 raw data was converted into subreads for subsequent analysis using the RS_Subreads.1 pipeline
427 within SMRTPortal (version 2.3), the web-based bioinformatics suite for analysis of RSII data.

428

429 **DNA extraction for PacBio (Sequel II) Sequencing (Data Source #2)**

430 For all samples from Azerbaijan and Georgia, MTB cultures were grown in 7H9+ADST broth to
431 A600 0.5–1.0. Pelleted cells were heat killed at 80°C for 2 hours. Cell pellets were resuspended in
432 450ul TE-Glu, 50ul of 10 mg/mL lysozyme was added and incubated at 37°C overnight. To each
433 sample 100ul of 10% sodium dodecyl sulfate and 50ul of 10 mg/ml proteinase K was added and
434 incubated at 55°C for 30 minutes. 200 ul of 5M sodium chloride and 160 ul Cetramide Saline
435 Solution (preheated 65°C) was added then incubated for 65°C for 10 minutes. To each sample 1
436 ml chloroform:isoamyl alcohol (24:1) was added, mixed gently by inversion. Samples were
437 centrifuged at 5000g for minutes, and 900ul of aqueous layer was transferred to fresh tube. DNA
438 was re-extracted with chloroform:isoamyl alcohol (24:1) and 800 ul of aqueous layer was
439 transferred to fresh tube. To 800 aqueous layer 560 ul isopropanol was added, mix gently by
440 inversion. The precipitated DNA was collected by centrifuging for 10 minutes and supernatant
441 was removed. DNA was washed with 70% ethanol, and DNA was collected by centrifuging and
442 supernatant removed. Air dried DNA pellet was dissolved overnight in 100 ul of TE buffer, and
443 stored at 4°C.

444
445 For all samples from Moldova, DNA was extracted according to CTAB protocol³².

446

447 **PacBio (Sequel II) Sequencing (Data Source #2)**

448 Approximately 1 µg of high molecular weight genomic DNA was used as input for SMRTbell
449 preparation according to the manufacturer's protocol (Preparing Multiplexed Microbial Libraries
450 Using SMRTbell Express Template Prep Kit 2.0, Pacific Biosciences). Briefly, HMW gDNA was
451 sheared to ~15kb using the Covaris g-tube at 2029 x g. For about half of the samples the
452 molecular weight of the DNA did not need shearing. Following shearing, gDNA underwent DNA
453 damage repair, ligation to SMRTbell barcoded adaptors and exonuclease treatment to remove
454 any unligated gDNA. At least 500 ng of pooled SMRTbell library per sample was cleaned with
455 AMPure PB beads and 7-50 kb fragments were size selected using the BluePippin system on 0.75%
456 agarose cassettes and S1 ladder, as specified by the manufacturer (Sage Science). The pool of
457 size-selected SMRTbell libraries were annealed to v4 sequencing primer and bound to the
458 polymerase prior to loading on the Sequel II sequencing system (Pacific Biosciences). Sequencing
459 was performed using version 1 chemistry and 15-hour movies.

460

461 **H37Rv reference genome and gene annotations**

462 The H37Rv (NCBI Accession: NC_000962.3) genome sequence and annotations was used as the
463 standard reference genome for all analyses. Functional category annotations for all genes of
464 H37Rv were downloaded from Release 3 (2018-06-05) of MycoBrowser³³
465 (<https://mycobrowser.epfl.ch/releases>). PE/PPE sub-family annotations of H37Rv were taken from

466 Ates et al.³⁴. Programmatic visualization of data along with annotations of the H37Rv genome
467 were made using the DNA Features Viewer python library³⁵.

468

469 **Genome assembly with PacBio long-read data**

470 All PacBio reads were assembled using Flye³⁶ (v2.6). After assembly, Flye performed three rounds
471 of iterative polishing of the genome assembly with the PacBio subreads, producing a polished de
472 novo PacBio assembly. If Flye identified the presence of a complete circular contig, Circlator³⁷
473 (v1.5.5) was used to standardize the start each assembly at the DnaA (Rv0001) locus.

474

475 **Polishing of *de novo* PacBio assemblies with Illumina WGS**

476 The paired-end Illumina WGS reads were trimmed with Trimmomatic³⁸ (v0.39) with the following
477 parameters: 2:30:10:2:true SLIDINGWINDOW:4:20 MINLEN:75. Trimmed reads were aligned to the
478 associated de novo PacBio assembly with BWA-MEM³⁹ (v0.7.17). Duplicate reads were removed
479 from the resulting alignments using PICARD⁴⁰ (v2.22.5). Using the deduplicated alignments, Pilon²²
480 (v1.23) was then used to correct SNPs and small INDELS in the *de novo* PacBio assembly, producing
481 a high confidence assembly polished by both PacBio and Illumina WGS.

482

483 **Identifying mixed infections using F2 metric and removing mismatched PacBio and** 484 **Illumina WGS**

485 To further reduce the effects of contamination, we used the F2 metric to identify samples that
486 may have inter-lineage variation due to co-infection⁴¹. The F2 metric measures the heterogeneity
487 of genotypes at known lineage defining positions of the H37Rv genome. We computed the F2
488 lineage-mixture metric for both PacBio and Illumina WGS from each isolate. Isolates were filtered
489 out if either the F2 metric for Illumina sequencing passed 0.05 or the F2 metric for PacBio
490 sequencing passed 0.35. The threshold used for PacBio sequencing subreads is much higher
491 because the inherent error rate per read is much higher than Illumina.

492

493 During polishing we identified the N0052 isolate from Chiner-Oms et al.²⁰ as a potential sample
494 mismatch, meaning PacBio and Illumina WGS were not performed on the same clinical isolate.
495 When polishing the de novo assembly of N0052, we found that the following changes were
496 performed based on the Illumina WGS: 594 SNPs, 19 insertions, and 92 deletions. The extreme
497 number of corrected SNPs by Illumina polishing is drastically different from the known error
498 profile (**Additional File 2-3**). Additionally, the inferred sub-lineage of the de novo PacBio
499 assembly was lineage 2.2.1, while the inferred sub-lineage based on Illumina WGS and the Illumina
500 Polished PacBio assembly was lineage 2.2.2 (**Additional File 2**). The fact that the polishing with
501 Illumina WGS changed known lineage defining SNPs makes the sample further suspect as a
502 mismatch. Thus, N0052 was removed from analysis as to minimize chances of benchmarking
503 wrongly matched data.

504

505 **Evaluation of PacBio genome assembly characteristics and multiple genome** 506 **alignment**

507 FastANI⁴² was used to calculate the average nucleotide identity to the H37Rv reference genome
508 for all completed genome assemblies. The Prokka (v1.13) genome annotation pipeline⁴³ was used
509 to annotate genes in each completed genome assembly. The genome size and GC content of the
510 entire genome was calculated from each assembly using custom python code. The
511 progressiveMauve algorithm of the Mauve (v2.4.0)⁴⁴ alignment software was used to perform
512 multiple sequence alignment of all 36 completed Mtb assemblies and the H37Rv reference
513 genome (NCBI Accession: NC_000962.3). The multiple genome alignments of H37Rv and 36
514 assemblies were visualized using the Mauve GUI⁴⁵ (**Supp. Figure 2**).

515

516 **Variant calling and structural variant detection using complete PacBio assemblies**

517 Minimap2⁴⁶ was used to align each polished circular completed assembly to the H37Rv reference
518 genome, producing a base-level alignment of similar regions of the assembly to H37Rv. In regions
519 with high sequence diversity or large structural variation, Minimap2 will not produce alignments.
520 To account for this, the NucDiff⁴⁷ analysis pipeline, which uses the MUMmer⁴⁸ aligner internally,
521 was also used to detect and classify the presence of large structural variants relative to the H37Rv
522 reference. All structural variants (≥ 50 bp) identified by NucDiff for each genome assembly can be
523 found in (**Additional File 16**).

524

525 **Illumina WGS data processing for variant calling relative to H37Rv**

526 Paired-end Illumina reads were trimmed with Trimmomatic (v0.39) with the following parameters:
527 2:30:10:2:true SLIDINGWINDOW:4:20 MINLEN:75. Trimmed reads were aligned to the H37Rv
528 reference genome (NC_000962.3) with BWA-MEM³⁹ (v0.7.17). Duplicate reads were removed from
529 the resulting alignments using PICARD⁴⁰ (v2.22.5). Using the deduplicated alignments, small
530 genome variants (SNSs and INDELS) were inferred using Pilon²² (v1.23). Samtools, Bcftools, and
531 BEDtools were used as needed for SAM/BAM, and VCF/BCF format file manipulation⁴⁹⁻⁵¹.

532

533 **Phylogenetic inference using complete genome assemblies**

534 All single nucleotide variants inferred through alignment with Minimap2 of PacBio assembly to
535 the H37Rv genome were concatenated across the 36 strains. Any SNS position which was ever
536 ambiguously called in at least 1 isolate was excluded (No NAs allowed, only REF or ALT alleles
537 allowed). Thus, in order for a SNS position to be included it needed to have no ambiguity relative
538 to the H37Rv reference in any isolate. FastTree⁵² was used to infer an approximate maximum
539 likelihood phylogeny from the concatenated SNS alignment of all 36 clinical Mtb isolates (15,673
540 total positions across 36 Mtb clinical isolates).

541

542 **Measuring repetitive sequence content of the H37Rv reference genome using Pileup** 543 **Mappability**

544 We evaluated sequence uniqueness using a *mappability* metric defined as the inverse of the
545 number of times a sequence of length K appears in a genome allowing for e mismatches and
546 considering the reverse complement⁵³. The *pileup mappability* of a position in a genome is then
547 defined as the average mappability of all overlapping k-mers. Thus, there are 2 parameters when
548 calculating mappability, k (length of k-mer) and e (number of base mismatches allowed in
549 counting matching k-mers). Genmap⁵⁴ (v1.3) was used to calculate the mappability of all k-mers
550 across the H37Rv reference genome with the following parameters: k-mer sizes of 50, 75, 100,
551 125, 150 base pairs and $E = 0-4$ mismatches. The Gene-level mappability ($k = 50$ bp , $e = 4$
552 mismatches) scores were computed as the average pileup mappability across all genes bodies
553 annotated in H37Rv (NCBI Accession: NC_000962.3). The base level pileup mappability scores of
554 H37Rv are available in TSV and BEDGRAPH format for easy visualization in a genome browser
555 **(Additional Files 6 and 17)**.

556

557 **Calculation of Empirical Base-level Recall (EBR) of Illumina variant calling**

558 The goal of the empirical base-level recall (EBR) for score is to summarize the consistency by which
559 Illumina WGS correctly evaluated any given genomic position. The EBR for a genomic position
560 was defined as the proportion isolates where Illumina WGS confidently and correctly agreed with
561 the PacBio defined ground truth. The ground truth was inferred for each isolate by directly
562 comparing the completed PacBio genome assembly to the H37Rv reference using Minimap2⁴⁶
563 and NucDiff⁴⁷. Due to Minimap2's inability to classify large structural variants, the ground truth
564 relative to H37Rv was supplemented with the structural variant calls generated by the NucDiff
565 analysis pipeline. Illumina WGS reads were aligned to the H37Rv reference genome with BWA-
566 MEM³⁹, and variants were inferred with the Pilon²² variant detection tool. In addition to identifying
567 variants relative to the reference genome, Pilon provides variant calling annotations for all
568 positions of H37Rv. The variant calling quality annotations of Pilon for all positions of H37Rv were
569 parsed for comparison to the PacBio defined ground truth for each isolate evaluated.

570 Only the following comparison outcomes were classified as a correctly recalled position:

571 1) Both Illumina variant calling and the PacBio ground truth agree on the genotype of a genomic
572 position, 2) Both Illumina variant calling and the PacBio ground truth agree that a genomic
573 position is deleted.

574

575 The following comparison outcomes were classified as poorly recalled position:

576 3) The PacBio ground truth supports a deletion, but Illumina is not confident in the presence of
577 the deletion, 4) Both Illumina variant calling and the PacBio ground truth disagree on the genotype
578 of a genomic position, 5) The PacBio ground truth supports the presence of a genomic region,
579 while Illumina variant calling did not confidently support the presence of the region. 6) Illumina

580 variant calling erroneously supports a deletion at a genomic position which is not deleted in the
581 PacBio ground truth.

582

583 The following EBR comparison outcomes were classified as ambiguous (N/A) due to ambiguities
584 in the interpretation of the ground truth: a) Cases where the PacBio ground truth contained
585 genome duplications relative to H37Rv, b) Cases where the PacBio ground truth did not provide
586 a confident alignment or structural variant call due to high sequence divergence from the
587 reference sequence.

588

589 For calculating the EBR for a genomic position, ambiguous (N/A) outcomes were ignored when
590 the number of N/As was $\leq 25\%$. In the case that a position had greater than 25% N/As at a
591 genomic position, the EBR score was defined as "Ambiguous". Ambiguous (N/A) EBR scores
592 represent locations of the H37Rv genome where there appeared to be systematic trouble in
593 determining the ground truth genotype.

594

595 The base level EBR scores are available in TSV and BEDGRAPH format for easy visualization in a
596 genome browser (**Additional Files 6 and 18**).

597

598 **Evaluating characteristics of low empirical performance across Mtb genome**

599 The Illumina WGS variant caller used, Pilon, produces VCF tags for all reference positions
600 evaluated, including positions which were confidently called a reference. The tags associated with
601 each position can either be PASS or a combination of non-pass tags (LowCov, Del, Amb). Each
602 genomic position can be assigned a combination of the following VCF Tags: a) PASS, signifying
603 confirmation of either a reference or an alternative allele. b) LowCov, signifying insufficient high
604 confidence reads (Depth < 5). c) Del, signifying that the position is confidently inferred to be
605 deleted. d) Amb, signifying evidence for more than one allele at this position. We quantified the
606 frequency of all combinations of these tags across all positions that were classified as "poor
607 recalled" during EBR evaluation.

608

609 **Measuring sequencing bias with per-base relative depth**

610 We measured sequencing bias using the relative depth statistic, which for a given genome
611 assembly and sequencing dataset, is defined as the sequencing depth per site divided by average
612 depth across the entire genome⁴. We evaluated the relative depth of all base pair positions of all
613 sequencing runs (Illumina and PacBio) relative to the corresponding isolates' complete PacBio
614 genome assembly. The sequencing depth of a base pair position was defined as the number of
615 reads with a nucleotide aligning to the position of interest. We calculated the mean coverage
616 across a sample by simply averaging the depth across all positions of the evaluated genome. For
617 ambiguous mapping reads, the aligners used (BWA-mem and Minimap2) use a random

618 assignment policy between all possible alignment locations. This allows for approximation of
619 depth in regions with non-uniquely mapping reads. For each individual Mtb isolate, we then
620 calculated the mean relative depth across all positions with the same GC content (100 bp window
621 size, **Additional File 8**).

622

623 **Defining and excluding ambiguous regions relative to H37Rv (per isolate genome** 624 **assembly)**

625 Following GA4GH (Global Alliance for Genomics & Health) benchmarking guidelines²³, we
626 excluded regions of the genome, where definition of the ground truth had ambiguity in its
627 definition relative to the reference genome. The following comparison outcomes were classified
628 as ambiguous (N/A) due to ambiguities in the interpretation of the ground truth: a) Cases where
629 the PacBio ground truth contained duplications relative to H37Rv, b) Cases where the PacBio
630 ground truth did not provide a confident alignment or structural variant call due to high sequence
631 divergence relative to H37Rv. These regions thus represent sequences of divergence relative to
632 the reference genome.

633

634 The percentage of the reference genome that was identified as “ambiguous” was consistently less
635 than 1% for all 36 clinical isolates evaluated. The median percent of the genome where the ground
636 truth was “ambiguously defined” was 0.4% (IQR: 0.3% - 0.5%). A large majority of these ambiguous
637 ground truth regions were either in Mobile Genetic Elements, PE_PGRS or PPE_MPTR genes. The
638 ambiguously defined regions for each isolate can be found in **Additional File 4**. Additionally, all
639 regions of the H37Rv genome which were ambiguous in over 25% of isolates, signifying high
640 levels of ambiguity, are present in **Additional File 5**.

641

642 **Defining the putative low confidence (PLC) regions of the H37Rv genome**

643 The regions most commonly excluded from Mtb genomics analysis, also referred to as the Putative
644 Low Confidence (PLC) regions in this work, were based on current literature^{16,24,55,56}. Specifically,
645 we defined the PLC regions as the union of the 168 PE/PPE genes, all mobile genetic elements
646 (MGEs), and 82 genes with repetitive content previously identified²⁴. PLC regions are defined in
647 **Additional File 19** (BED format). Non-PLC regions were simply defined as the complement of the
648 PLC genes.

649

650 **Evaluating variant calling performance of genome masking approaches**

651 Following the small variant benchmarking standards outlined by the GA4GH, we used Hap.py
652 (v0.3.13) to evaluate the Illumina WGS variant calling performance of Pilon for all 36 isolates
653 individually. For each complete genome assembly, SNSs and small INDELS 1-15 bp inferred by the
654 Minimap2-paftools pipeline were used as ground truth. We evaluated variant calling performance
655 of Illumina WGS when using different region filtering schemas: (1) masking of all PLC genes, the

656 current standard practice, (2) masking of repetitive regions with P-Map-K50E4 < 100%, and (3) No
657 masking. Masking schemas (1 and 2) are provided in BED format (**Additional File 19 and 20**).
658 After applying each masking schema, we filtered potential variants according to whether the Pilon
659 variant calling pipeline gave the variant a PASS filter and the mean mapping quality (MQ) of all
660 reads aligned to the variant position.

661
662 For each combination of region masking and variant filtering using mapping quality, we then
663 calculated the absolute number of true positives (TP, i.e. a variant in the ground truth variant set
664 and correctly called by the Illumina variant calling pipeline), false positives (FP, the Illumina variant
665 calling pipeline calls a variant not in the ground truth set), and false negative (FN, the variant is in
666 the ground truth set but is not called by the Illumina variant calling pipeline) variant calls. For each
667 set of parameters, we calculated the overall precision (positive predictive value) as $TP/(TP + FP)$,
668 and recall (sensitivity) as $TP/(TP + FN)$. In agreement with the default behavior of Hap.py, and to
669 avoid undefined precision values, filtering parameters that yielded no TP or FP were defined as
670 having a precision of 1.0 and a recall of 0. Additionally, we calculated the F1-score, which weights
671 precision and recall with equally: $F1 = 2 * (precision * recall)/(precision + recall)$. The F1 score
672 summarizes each variant calling performance as a single value between 0 and 1 (where 1
673 represents both perfect precision and recall).

674
675 To aggregate the performance evaluation across all 36 isolates, the mean and standard error of
676 the mean (SEM) of precision, recall and F1 score was calculated for all sets of parameters evaluated
677 (**Additional File 10**). The individual variant calling performance statistics for each isolate can also
678 be found in **Additional File 10**. The variant calling performance comparison of shorter (1-5bp) vs
679 longer (6-15bp) INDELs can be found in **Additional File 11**.

680
681 **Evaluating variant calling performance near regions with structural variation and**
682 **repetitive sequence content**

683 Using Hap.py and the same approach defined in the above section, we evaluated SNS variant
684 calling performance in the following types of regions: (1) SNSs in regions with perfect mappability
685 (Pmap-K50E4 = 1) with no identified SV (2) SNSs in regions with low mappability (Pmap-K50E4 <
686 1) with no identified SV, (3) SNSs in regions with perfect mappability within 100 bp of any
687 identified SV, and (4) SNSs in regions with low mappability within 100bp of any identified SV.
688 Genomic contexts not near SVs (1 and 2) were evaluated with MQ thresholds ranging from 1-60.
689 For genomic contexts within 100 bp of an SV (3 and 4), the MQ thresholds evaluated ranged from
690 1-40. The MQ threshold evaluated near SVs was limited due to the fact that a majority of SNSs
691 near SVs typically have lower MQ values, and higher MQ values resulted in recalls of approximately
692 0. As explained in the previous section, the mean and SEM of precision, recall, and F1 score were
693 calculated for all MQ filtering thresholds across all 4 region types (**Additional File 12**).

694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728

Evaluation of the distribution of potential false positive SNS calls across the Mtb genome

False positive SNS calls were identified by the Hap.py evaluation software through comparison to the assembly-based ground truth variant call set. Additionally, false positive calls with $MQ < 30$ were filtered out, as to only include false positives which would realistically pass standard filtering. For each genomic region (gene or intergenic region) of the H37Rv genome, the total number of overlapping false positives across all 36 isolates was calculated (**Additional File 9**). Across all 36 clinical isolates, there were 548 false positive SNSs with $MQ \geq 30$ and 696 total false positive SNS with $MQ \geq 1$ detected.

Defining Refined Low Confidence (RLC) regions

We defined the refined low confidence regions (RLC) of the Mtb reference genome as the union of A) The 30 false positive hot spot regions (gene and intergenic) identified (65 kb), B) poorly recalled genomic regions as identified by EBR ($EBR < 0.9$, 142 kb), and C) regions with frequently ambiguously defined ground truths (16 kb). We provide the complete set of RLC regions in BED format (177 kb, **Additional File 13**), along with each separate component of the RLC regions in BED format (**Additional Files 21, 22, and 23**). For very conservative masking of the Mtb reference genome, we additionally provide a masking scheme that specifies the union of a) the RLC regions and b) all low pileup mappability regions ($PmapK50E4 < 1$) (277 kb, **Additional File 14**).

729 **References**

- 730 1. Li, H., Ruan, J. & Durbin, R. Mapping short DNA sequencing reads and calling variants using
731 mapping quality scores. *Genome Res.* **18**, 1851–1858 (2008).
- 732 2. Li, H. Toward better understanding of artifacts in variant calling from high-coverage
733 samples. *Bioinformatics* **30**, 2843–2851 (2014).
- 734 3. Nakamura, K. *et al.* Sequence-specific error profile of Illumina sequencers. *Nucleic Acids Res.*
735 **39**, e90 (2011).
- 736 4. Ross, M. G. *et al.* Characterizing and measuring bias in sequence data. *Genome Biol.* **14**, R51
737 (2013).
- 738 5. Goig, G. A., Blanco, S., Garcia-Basteiro, A. L. & Comas, I. Contaminant DNA in bacterial
739 sequencing experiments is a major source of false genetic variability. *BMC Biol.* **18**, 24
740 (2020).
- 741 6. Barbitoff, Y. A. *et al.* Systematic dissection of biases in whole-exome and whole-genome
742 sequencing reveals major determinants of coding sequence coverage. *Sci. Rep.* **10**, 2057
743 (2020).
- 744 7. Modlin, S. J. *et al.* Exact mapping of Illumina blind spots in the *Mycobacterium tuberculosis*
745 genome reveals platform-wide and workflow-specific biases. *Microb Genom* (2021)
746 doi:10.1099/mgen.0.000465.
- 747 8. Benjamini, Y. & Speed, T. P. Summarizing and correcting the GC content bias in high-
748 throughput sequencing. *Nucleic Acids Res.* **40**, e72 (2012).
- 749 9. Aird, D. *et al.* Analyzing and minimizing PCR amplification bias in Illumina sequencing
750 libraries. *Genome Biol.* **12**, R18 (2011).
- 751 10. Paten, B., Novak, A. M., Eizenga, J. M. & Garrison, E. Genome graphs and the evolution of
752 genome inference. *Genome Res.* **27**, 665–676 (2017).
- 753 11. Garrison, E. *et al.* Variation graph toolkit improves read mapping by representing genetic
754 variation in the reference. *Nat. Biotechnol.* **36**, 875–879 (2018).
- 755 12. Schmid, M. *et al.* Pushing the limits of de novo genome assembly for complex prokaryotic
756 genomes harboring very long, near identical repeats. *Nucleic Acids Res.* **46**, 8953–8965
757 (2018).
- 758 13. Rhoads, A. & Au, K. F. PacBio Sequencing and Its Applications. *Genomics Proteomics*
759 *Bioinformatics* **13**, 278–289 (2015).
- 760 14. Wenger, A. M. *et al.* Accurate circular consensus long-read sequencing improves variant
761 detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
- 762 15. De Maio, N. *et al.* Comparison of long-read sequencing technologies in the hybrid assembly
763 of complex bacterial genomes. *Microb Genom* **5**, (2019).
- 764 16. Meehan, C. J. *et al.* Whole genome sequencing of *Mycobacterium tuberculosis*: current
765 standards and open issues. *Nat. Rev. Microbiol.* (2019) doi:10.1038/s41579-019-0214-5.

- 766 17. Coscolla, M. & Gagneux, S. Consequences of genomic diversity in *Mycobacterium*
767 tuberculosis. *Semin. Immunol.* **26**, 431–444 (2014).
- 768 18. Hicks, N. D. *et al.* Clinically prevalent mutations in *Mycobacterium tuberculosis* alter
769 propionate metabolism and mediate multidrug tolerance. *Nat Microbiol* **3**, 1032–1042
770 (2018).
- 771 19. Holt, K. E. *et al.* Frequent transmission of the *Mycobacterium tuberculosis* Beijing lineage
772 and positive selection for the EsxW Beijing variant in Vietnam. *Nat. Genet.* **50**, 849–856
773 (2018).
- 774 20. Chiner-Oms, Á. *et al.* Genome-wide mutational biases fuel transcriptional diversity in the
775 *Mycobacterium tuberculosis* complex. *Nat. Commun.* **10**, 3994 (2019).
- 776 21. Ngabonziza, J. C. S. *et al.* A sister lineage of the *Mycobacterium tuberculosis* complex
777 discovered in the African Great Lakes region. *Nat. Commun.* **11**, 2917 (2020).
- 778 22. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
779 and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 780 23. Krusche, P. *et al.* Best practices for benchmarking germline small-variant calls in human
781 genomes. *Nat. Biotechnol.* **37**, 555–560 (2019).
- 782 24. Coscolla, M. *et al.* M. tuberculosis T Cell Epitope Analysis Reveals Paucity of Antigenic
783 Variation and Identifies Rare Variable TB Antigens. *Cell Host Microbe* **18**, 538–548 (2015).
- 784 25. Thomas, S. K. *et al.* Modern and ancestral genotypes of *Mycobacterium tuberculosis* from
785 Andhra Pradesh, India. *PLoS One* **6**, e27584 (2011).
- 786 26. Sharifipour, E., Farnia, P., Mozafari, M., Irani, S. & Akbar Velayati, A. Deletion of region of
787 difference 181 in *Mycobacterium tuberculosis* Beijing strains. *Int J Mycobacteriol* **5 Suppl 1**,
788 S238–S239 (2016).
- 789 27. Walter, K. S. *et al.* Genomic variant-identification methods may alter *Mycobacterium*
790 tuberculosis transmission inferences. *Microb Genom* **6**, (2020).
- 791 28. Jajou, R. *et al.* Towards standardisation: comparison of five whole genome sequencing
792 (WGS) analysis pipelines for detection of epidemiologically linked tuberculosis cases. *Euro*
793 *Surveill.* **24**, (2019).
- 794 29. Farhat, M. R. *et al.* GWAS for quantitative resistance phenotypes in *Mycobacterium*
795 tuberculosis reveals resistance genes and regulatory regions. *Nat. Commun.* **10**, 2128
796 (2019).
- 797 30. Rosenthal, A. *et al.* The TB Portals: an Open-Access, Web-Based Platform for Global Drug-
798 Resistant-Tuberculosis Data Sharing and Analysis. *J. Clin. Microbiol.* **55**, 3267–3282 (2017).
- 799 31. Epperson, L. E. & Strong, M. A scalable, efficient, and safe method to prepare high quality
800 DNA from mycobacteria and other challenging cells. *J Clin Tuberc Other Mycobact Dis* **19**,
801 100150 (2020).
- 802 32. Wilson, K. Preparation of genomic DNA from bacteria. *Curr. Protoc. Mol. Biol.* **Chapter 2**,
803 Unit 2.4 (2001).

- 804 33. Kapopoulou, A., Lew, J. M. & Cole, S. T. The MycoBrowser portal: a comprehensive and
805 manually annotated resource for mycobacterial genomes. *Tuberculosis* **91**, 8–13 (2011).
- 806 34. Ates, L. S. New insights into the mycobacterial PE and PPE proteins provide a framework for
807 future research. *Mol. Microbiol.* (2019) doi:10.1111/mmi.14409.
- 808 35. Zulkower, V. & Rosser, S. DNA Features Viewer: a sequence annotation formatting and
809 plotting library for Python. *Bioinformatics* **36**, 4350–4352 (2020).
- 810 36. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using
811 repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
- 812 37. Hunt, M. *et al.* Circlator: automated circularization of genome assemblies using long
813 sequencing reads. *Genome Biol.* **16**, 294 (2015).
- 814 38. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence
815 data. *Bioinformatics* **30**, 2114–2120 (2014).
- 816 39. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
817 *arXiv [q-bio.GN]* (2013).
- 818 40. Picard Tools - By Broad Institute. <http://broadinstitute.github.io/picard/>.
- 819 41. Wyllie, D. H. *et al.* Identifying Mixed Mycobacterium tuberculosis Infection and Laboratory
820 Cross-Contamination during Mycobacterial Sequencing Programs. *J. Clin. Microbiol.* **56**,
821 (2018).
- 822 42. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High throughput
823 ANI analysis of 90K prokaryotic genomes reveals clear species boundaries. *Nat. Commun.* **9**,
824 5114 (2018).
- 825 43. Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**, 2068–2069
826 (2014).
- 827 44. Darling, A. E., Mau, B. & Perna, N. T. progressiveMauve: multiple genome alignment with
828 gene gain, loss and rearrangement. *PLoS One* **5**, e11147 (2010).
- 829 45. Darling, A. C. E., Mau, B., Blattner, F. R. & Perna, N. T. Mauve: multiple alignment of
830 conserved genomic sequence with rearrangements. *Genome Res.* **14**, 1394–1403 (2004).
- 831 46. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100
832 (2018).
- 833 47. Khelik, K., Lagesen, K., Sandve, G. K., Rognes, T. & Nederbragt, A. J. NucDiff: in-depth
834 characterization and annotation of differences between two sets of DNA sequences. *BMC*
835 *Bioinformatics* **18**, 338 (2017).
- 836 48. Kurtz, S. *et al.* Versatile and open software for comparing large genomes. *Genome Biol.* **5**,
837 R12 (2004).
- 838 49. Danecsek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, (2021).
- 839 50. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and
840 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–
841 2993 (2011).

- 842 51. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
843 features. *Bioinformatics* **26**, 841–842 (2010).
- 844 52. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood trees
845 for large alignments. *PLoS One* **5**, e9490 (2010).
- 846 53. Derrien, T. *et al.* Fast computation and applications of genome mappability. *PLoS One* **7**,
847 e30377 (2012).
- 848 54. Pockrandt, C., Alzamel, M., Iliopoulos, C. S. & Reinert, K. GenMap: Ultra-fast Computation of
849 Genome Mappability. *Bioinformatics* (2020) doi:10.1093/bioinformatics/btaa222.
- 850 55. Lee, R. S., Proulx, J.-F., McIntosh, F., Behr, M. A. & Hanage, W. P. Previously undetected
851 super-spreading of Mycobacterium tuberculosis revealed by deep sequencing. *Elife* **9**,
852 e53245 (2020).
- 853 56. Coscolla, M. *et al.* Phylogenomics of Mycobacterium africanum reveals a new lineage and a
854 complex evolutionary history. *Microb Genom* (2021) doi:10.1099/mgen.0.000477.
- 855 57. Borrell, S. *et al.* Reference set of Mycobacterium tuberculosis clinical strains: A tool for
856 research and product development. *PLoS One* **14**, e0214088 (2019).
- 857 58. Köster, J. & Rahmann, S. Snakemake--a scalable bioinformatics workflow engine.
858 *Bioinformatics* **28**, 2520–2522 (2012).
- 859 59. Marin, M. G. *Additional File 6 - Base level analysis of Empirical Base Pair Recall, Pileup*
860 *Mappability, and GC content across the H37Rv genome.* (2021). doi:10.5281/zenodo.4662193.

861
862

863 **Author Contributions**

864 MGM and MRF conceived, designed and conducted the study. MGM and MRF wrote the
865 manuscript with input from all authors. RVJ provided bioinformatics support and input on data
866 analysis. LEE, DD, M. Salfinger and M. Strong cultured Mtb isolates and performed DNA extraction
867 in preparation for PacBio sequencing of Dataset #1. IA, SV, and VC cultured Mtb isolates and
868 performed DNA extraction in preparation for PacBio sequencing of Dataset #2.. AR, MH, and BJ
869 selected clinical isolates and assisted in data processing for PacBio sequencing of Dataset #2. ZI
870 provided help and advice throughout the project. The final manuscript was read and approved by
871 all authors.

872 **Competing Interests**

873 The authors declare that they have no competing interests.

874 **Data availability and materials**

875 All new sequencing data generated for this study and complete Mtb genome assemblies were
876 submitted to NCBI SRA and Genbank databases under BioProject accession number PRJNA719670
877 (Submission Pending). The publicly available PacBio and Illumina data from two previously
878 published studies^{20,21,57} is available from PRJEB8783, PRJEB31443, PRJEB27802, and PRJNA598991.
879 SRA/ENA accessions and related sequencing metadata for all data can be found in Additional File
880 15. All code for data processing and analysis in this study is available from the following GitHub
881 repository, <https://github.com/farhat-lab/mtb-illumina-wgs-evaluation>. The repository README
882 provides instructions to run each part of the analysis using the Snakemake⁵⁸ workflow engine and
883 using Python based Jupyter notebooks.

884 **Acknowledgments**

885 We are grateful to Natalia Quiñones, and Karel Brinda for their helpful discussions and advise
886 throughout the project. We grateful to Melissa Smith and Irina Oussenko for their assistance in
887 PacBio (RS II) long read sequencing of the M. tuberculosis genomic DNA. We acknowledge NIH
888 Intramural Sequencing Center (NISC) for the PacBio (Sequel II) long-read sequencing of the M.
889 tuberculosis genomic DNA; Critical Path Institute (C-Path) and Translational Genomics Research
890 Institute (T-Gen) for the Illumina sequencing of the M. tuberculosis DNAs and for the mTB DNA
891 long-term storage; the International Science and Technology Center for their support in
892 establishing the TB Portal agreement with Georgia; CRDF Global for their support in establishing
893 the TB Portal agreements with Azerbaijan and Moldova. This research was supported in part by
894 the Office of Science Management and Operations of the NIAID.

895

896

897

898

899

900

901

902

903

904

905

906

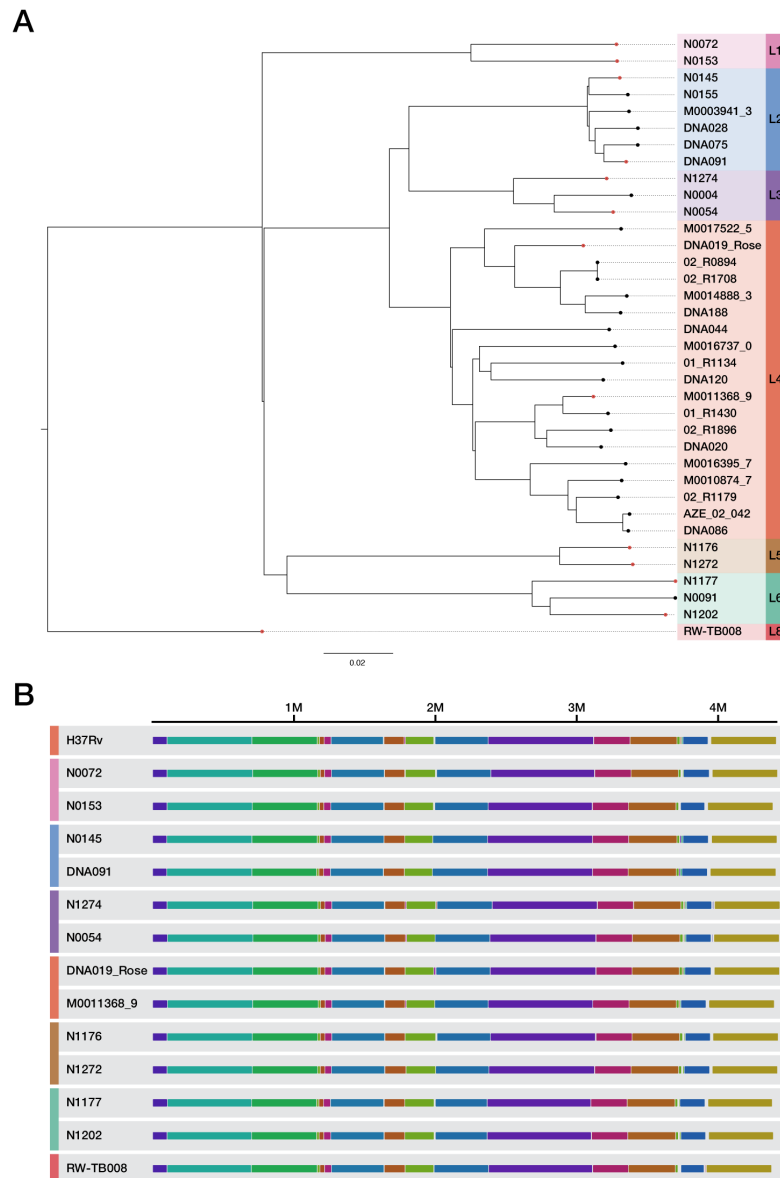
907

908

909

910 Figures & Tables

911 Figure 1

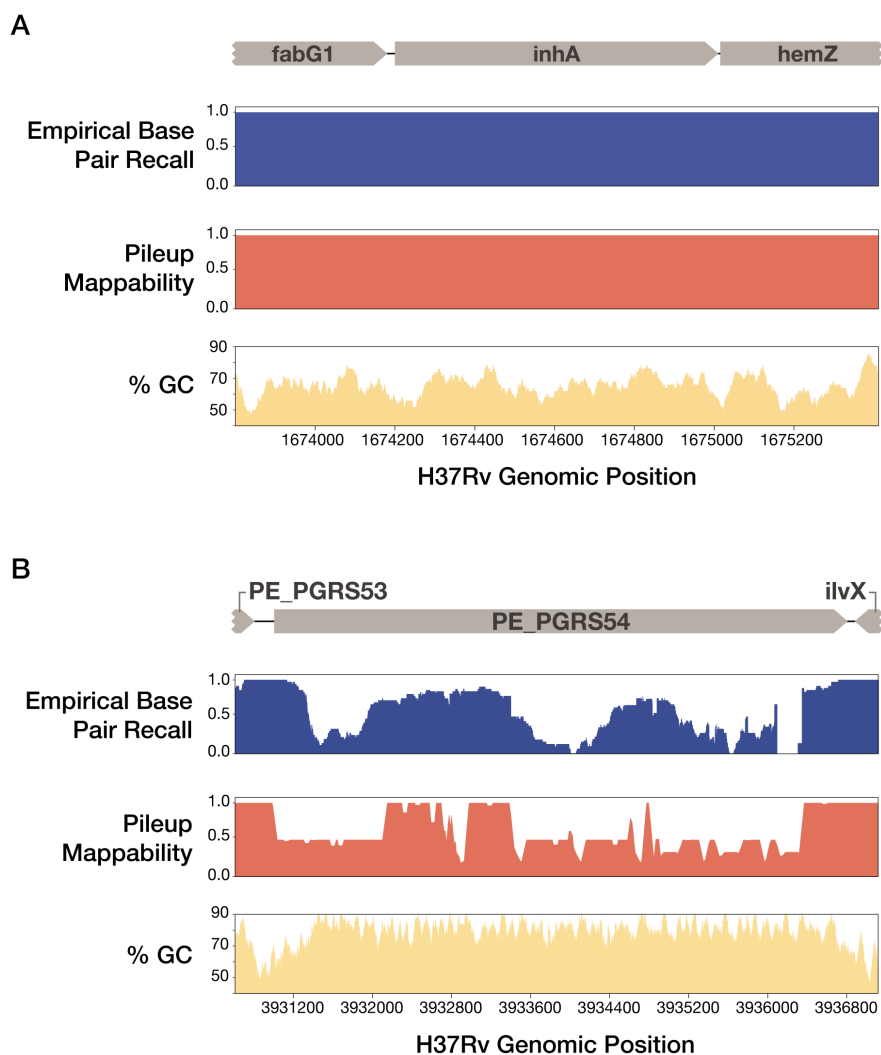


912

913 **Overview of 36 clinical *Mtb* isolates with completed genome assemblies. a)** Maximum
914 likelihood Phylogeny of *M. tuberculosis* isolates with PacBio complete genome assemblies. The
915 sequences of all 36 complete *M. tb* genomes were aligned to the H37rv reference genome using
916 minimap2, and a maximum likelihood phylogeny was inferred using a concatenated SNS alignment
917 (15,673 total positions). **b)** Representative isolates from each lineage sampled from the whole
918 genome sequence alignment between the H37Rv reference genome and all completed circular
919 *Mtb* genome assemblies, The complete alignment is visualized in Supplemental Figure 2. The
920 whole genome multiple sequence alignment was performed using the *progressiveMauve*⁴⁴
921 algorithm. Each contiguously colored region is a locally collinear block (LCB), a region without
922 rearrangement of homologous backbone sequence.

923
924

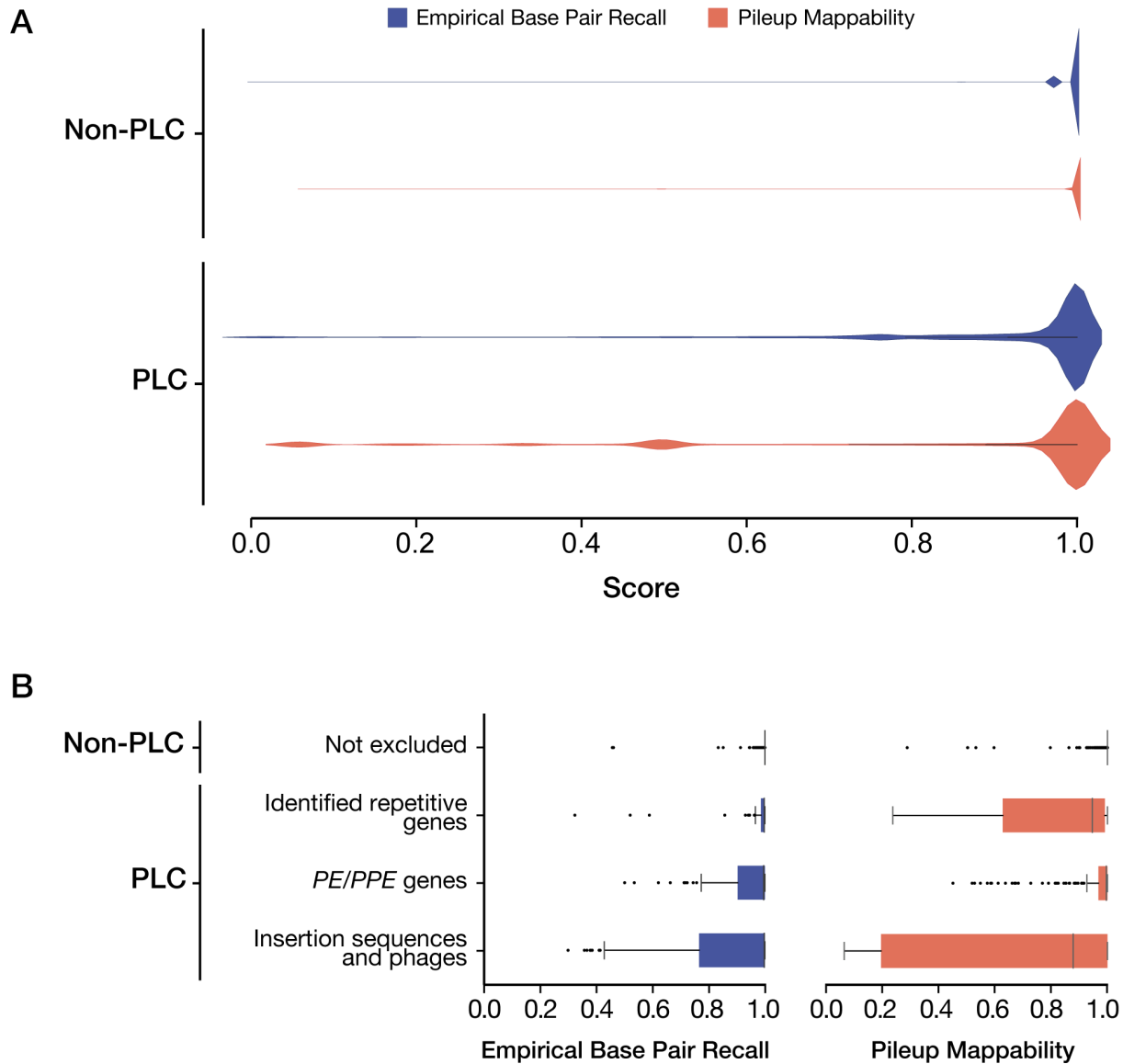
Figure 2



925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940

EBR, Pileup Mappability, and GC content across two example regions of the H37Rv genome. Empirical Base Pair Recall (EBR), Pileup Mappability (K=50 bp, e = 4 mismatches) and GC% (100 bp window) values are plotted across all base pair positions of two regions of interest. **a)** *InhA*, an antibiotic resistance gene, shows perfect EBR across the entire gene body. **b)** In contrast, *PE_PGRS54*, a known highly repetitive gene with high GC content, has extremely low EBR across the entire gene body. Browser tracks of EBR and Pileup Mappability in BEDGRAPH format are made available as Additional Files 17 and 18.

941 **Figure 3**



942 **The Distribution of EBR and Pileup Mappability scores in PLC and non-PLC regions.** **a)** The distribution of
 943 Empirical Base Pair Recall (EBR) and Pileup Mappability (P-Map, K=50,E=4) scores of PLC and non-PLC regions.
 944 Excluded regions harbor significantly more low EBR base pair positions when compared to the included genes,
 945 but 68% of routinely excluded positions still have $\geq 97\%$ EBR. The Pileup mappability with K=50 bp is lower in
 946 PLC regions (mean = 0.86) than non-PLC regions (mean = .997). **b)** The Distribution of gene-level mean EBR
 947 and P-Map (K=50,E=4) between PLC and non-PLC regions. We compared the mean EBR and Pileup
 948 Mappability across all genes within PLC and non-PLC regions. The *pe* and *ppe* gene families (PE/PPEs) and
 949 mobile genetic elements (MGE), which make up 82% of PLC genes, demonstrated significantly lower mean EBR
 950 and Pileup Mappability than other non-PLC genes.
 951

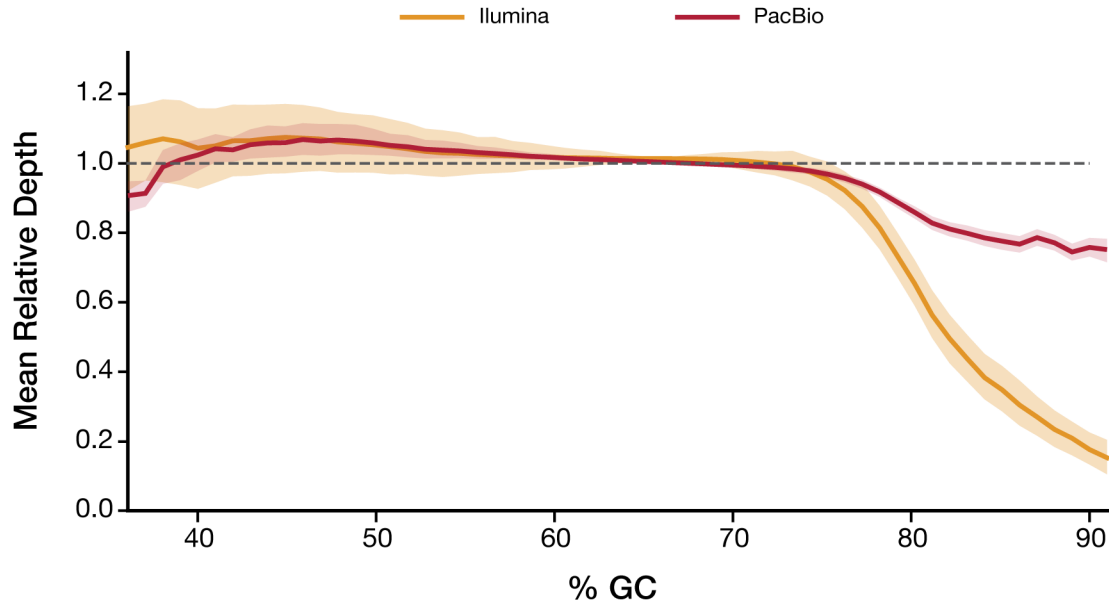
952

953

954

955

956 **Figure 4**



957 **Relative sequencing depth as a function of local GC content across all 36 complete isolates.** We evaluated the
958 relative depth of our Illumina and PacBio sequencing data as a function of GC content (100 bp window size) across all
959 positions of each isolate's complete genome assembly. The relative depth was averaged across all positions with the
960 same GC% across each genome assembly. The standard error of the mean of the relative depth across all 36 isolates is
961 shaded for each sequencing technology. At high (>70%) GC contents, Illumina starts to show lower relative depth
962 compared to PacBio sequencing.
963

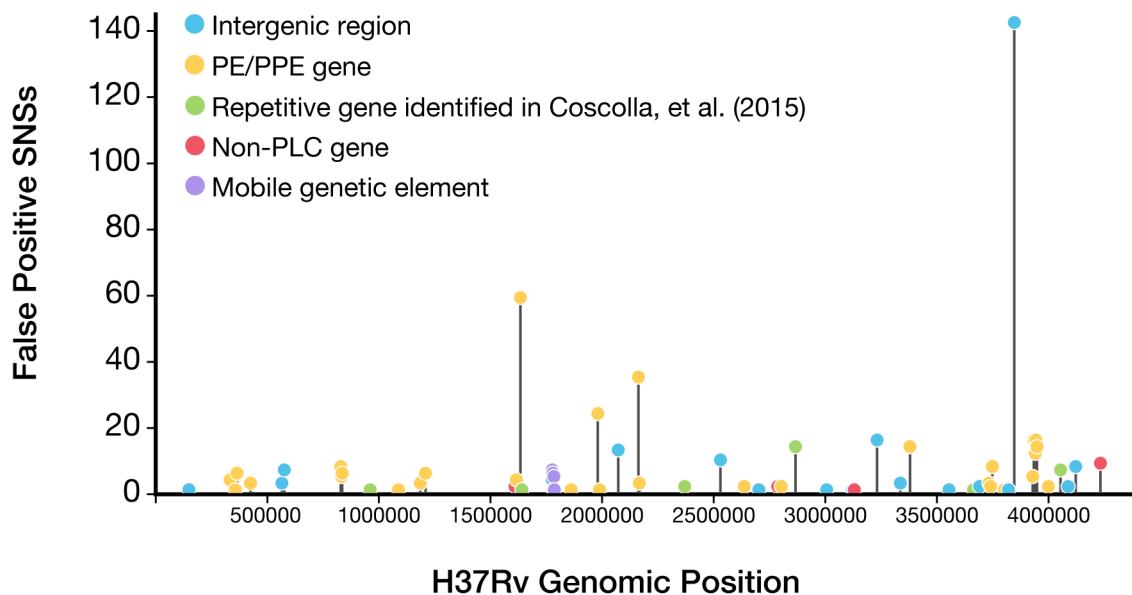
964

965

966

967

968 **Figure 5**

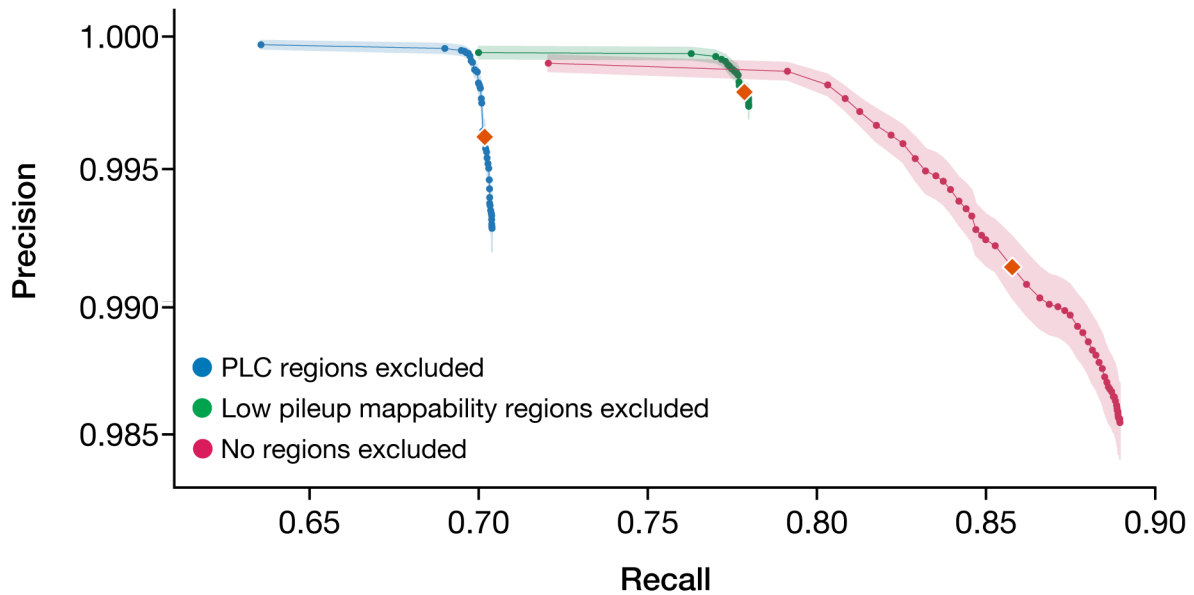


969 **The distribution of potential false positive SNS calls across all genomic regions of the H37Rv genome.** The
970 frequency of false positive SNS calls detected ($MQ \geq 30$) across all 36 isolates evaluated was plotted for all regions of
971 the H37Rv genome (gene or intergenic regions). The top 30 regions ranked by the number of total false positives
972 contained 89.4% (490/548) of the total false positive SNSs and spanned only 65 kb of the H37Rv genome. Full results
973 for all annotated genomic regions (gene or intergenic) can be found in Additional File 9.

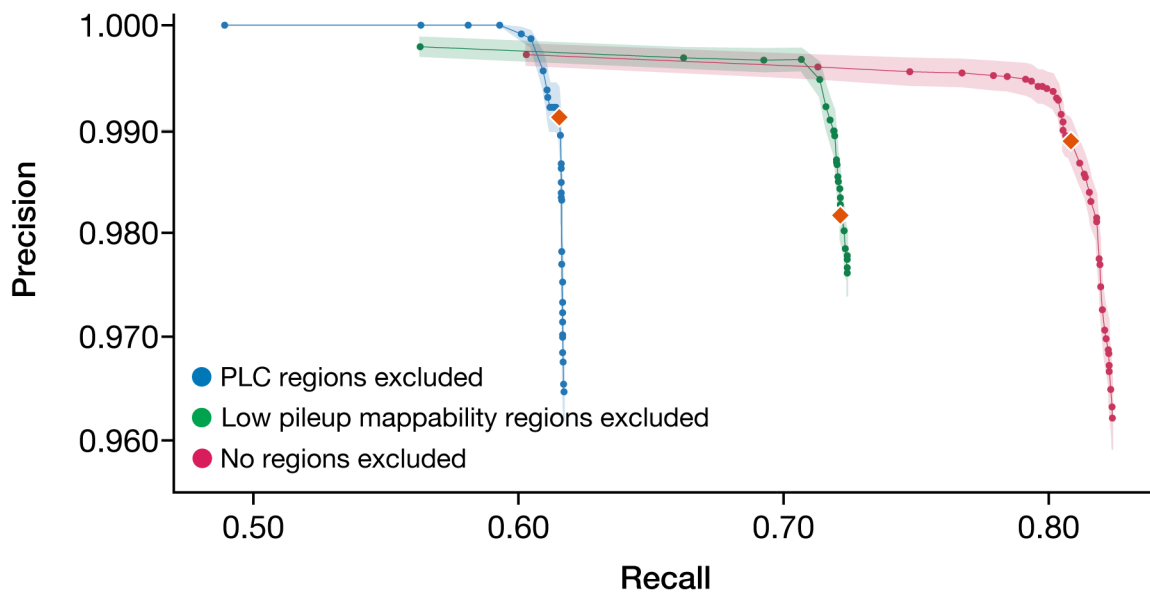
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990

991 **Figure 6**

A



B



992
993 **Mean SNV and INDEL variant calling performance across different masking approaches. a)** SNS variant calling
994 performance was evaluated across the following three schemas: (1) masking of regions with non-unique sequence, as
995 defined as positions with P-Map-K50E4 < 1, (2) No *a priori* masking of any regions, and compared to (3) masking of
996 all PLC genes (the current standard practice). **(b)** short INDEL (1-15 bp) variant calling performance was evaluated
997 across the same schemas. The orange diamonds represent the mean precision and recall using a MQ threshold of 40
998 for both (a) and (b). Shaded regions represent the SEM of precision across all 36 isolates evaluated.

999 For all masking approaches evaluated, the MQ thresholds evaluated ranged from 1-60. Complete benchmarking
1000 results can be found for each individual isolate in Additional File 10.

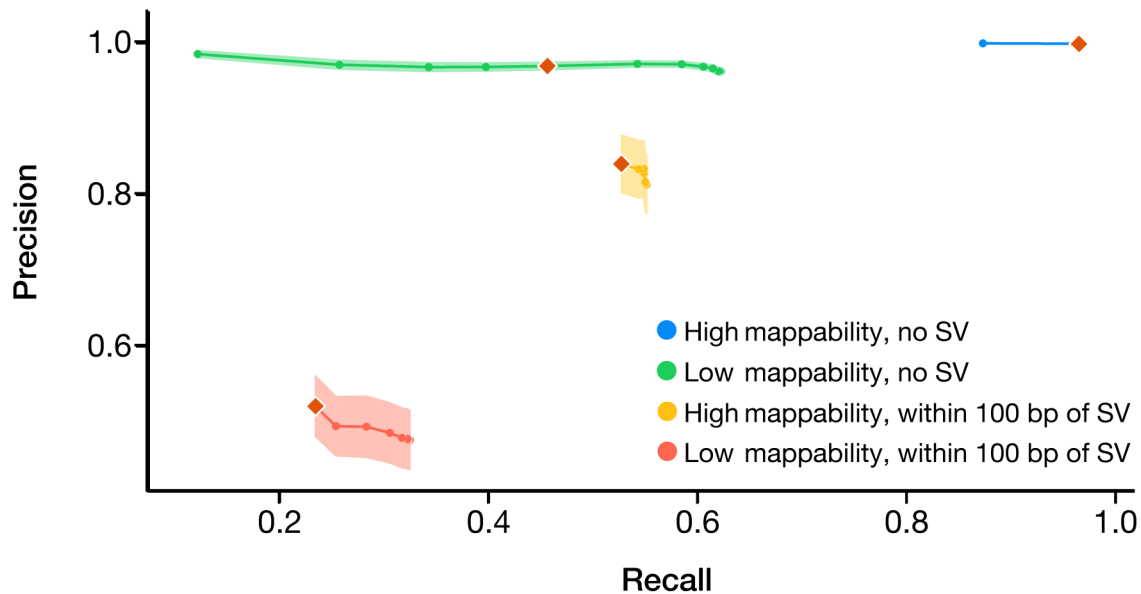
1001 **Table 1.**

Masking Schema	Metric Optimized	MQ Threshold	F1	Precision	Recall
Masking non-unique regions	F1-score	19	0.87	99.77%	77.98%
	Comparator	40	0.88	99.79%	77.86%
	Precision	60	0.82	99.94%	70.00%
No masking	F1-score	8	0.94	98.56%	88.95%
	Comparator	40	0.92	99.13%	85.77%
	Precision	60	0.83	99.90%	72.06%
Masking PLC genes (current standard)	F1-score	35	0.82	99.50%	70.30%
	Comparator	40	0.82	99.62%	70.17%
	Precision	60	0.77	99.97%	63.56%

1002 **Comparison of performance of proposed genome-masking schemas for SNS variant calling.** For each masking
 1003 scheme and MQ filtering threshold shown, the corresponding mean Precision, Recall, and F1 score is shown across all
 1004 36 Mtb isolates. Corresponding Precision-Recall curves are given in Figure 5A. Performance at a threshold of MQ \geq 40
 1005 is given as a common point of comparison across the three masking schemas.
 1006

1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021

1022 **Figure 7**



1023

1024

Variant calling performance of single nucleotide substitutions stratified by proximity to structural variants

1025

and low pileup mappability sequence. Mappability is dichotomized at Pmap-K50E4 =100% or <100%. Regions

1026

within 100bp of a SV categorized as "with SV". Precision and recall is plotted for the following genomic contexts: (1)

1027

regions with high mappability with no SV (Blue, F1 = 0.98 (precision = 99.89%, recall = 96.49%, MQ threshold of 40)),

1028

(2) regions with low mappability and no SV (green, F1 = 0.62 (precision = 96.98%, recall = 45.65%, MQ threshold of

1029

40), (3) regions with high mappability with SV (orange, F1 = 0.64 (precision = 84.07%, recall = 52.73%, MQ threshold

1030

of 40), (4) regions with low mappability and with SV (red, F1 = 0.32 (precision = 52.10%, recall = 23.47%, MQ

1031

threshold of 40). The standard error of the mean (SEM) for precision is shaded for each curve. Orange diamonds

1032

represent the precision and recall using the same MQ threshold of 40. Genomic contexts not near SVs (1 and 2) were

1033

evaluated with MQ thresholds ranging from 1-60. For genomic contexts within 100 bp of an SV (3 and 4), the MQ

1034

thresholds evaluated ranged from 1-40. Complete benchmarking results can be found for each individual isolate in

1035

Additional File 12.

1036

1037

1038

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048 **Supplementary Information**

1049 Additional File 1: Supplementary Figures and Tables (Figure S1-7, Table S1-6)

1050

1051 Additional File 2: Results and quality control for assembly and sequencing for both PacBio and Illumina
1052 sequencing

1053

1054 Additional File 3: List of all changes made during Illumina polishing of the *de novo* PacBio assemblies

1055

1056 Additional File 4: List of genomic regions with ambiguously defined ground truths relative to H37Rv for all
1057 each isolate evaluation

1058

1059 Additional File 5: List of genomic regions which were frequently had an ambiguously defined ground truth

1060

1061 Additional File 6: Table containing the EBR, Pileup Mappability, and GC% of all genomic positions of the
1062 H37Rv reference. Due to large file size, Additional File 6⁵⁹ is hosted on Zenodo at

1063 <https://zenodo.org/record/4662193>.

1064

1065 Additional File 7: EBR, and Pileup Mappability across all genomic regions of H37Rv (both genes and
1066 intergenic regions)

1067

1068 Additional File 8: Table of the mean relative sequencing depth of both Illumina and PacBio at varying GC%
1069 across all 36 isolates evaluated.

1070

1071 Additional File 9: Table containing the frequency of observed False Positive SNSs ($MQ \geq 30$) across all
1072 genomic regions of H37Rv (both genes and intergenic regions)

1073

1074 Additional File 10: Variant call benchmarking of SNSs and small indels (≤ 15 bp)

1075

1076 Additional File 11: Variant call benchmarking stratified by shorter (< 6 bp) and longer indels (6-15bp)

1077

1078 Additional File 12: Variant call benchmarking of SNSs stratified by proximity to an SV and low pileup
1079 mappability

1080

1081 Additional File 13: Masking scheme in BED format specifying the Refined Low Confidence Regions

1082

1083 Additional File 14: Masking scheme in BED format specifying the union of a) Refined Low Confidence
1084 Regions, and b) regions with Pileup Mappability ($K = 50$ bp, $E = 4$ mismatches) < 1 .

1085

1086

1087

1088

1089

1090 Additional File 15: SRA/ENA sequencing run metadata for PacBio and Illumina sequencing used in this
1091 study
1092
1093 Additional File 16: All identified structural variants for each complete genome assembly as identified by
1094 the NucDiff analysis pipeline.
1095
1096 Additional File 17: Base-level Pileup Mappability scores (P-Map-K50E4) across the H37Rv in BEDGRAPH
1097 format
1098
1099 Additional File 18: Base-level EBR scores (36 isolates) across the H37Rv in BEDGRAPH format
1100
1101 Additional File 19: Masking scheme for the Putative Low Confidence (PLC) Regions in BED format
1102
1103 Additional File 20: All regions with low pileup mappability (P-Map-K50E4 < 100%) in BED format
1104
1105 Additional File 21: Component (A) of RLC regions. Masking scheme Specifying the 30 false positive hot
1106 spot regions (gene and intergenic) in BED format.
1107
1108 Additional File 22:
1109 Component (B) of RLC regions. Masking scheme specifying poorly recalled genomic regions as identified
1110 by EBR < 0.9) in BED format.
1111
1112 Additional File 23:
1113 Component (C) of RLC regions. Masking scheme specifying regions that frequently (> 25%) had an
1114 ambiguously defined ground truth in BED format. Same information as Additional File 5 but this file is
1115 instead in BED format.