

1 **Performance characteristics of next-generation sequencing for antimicrobial resistance gene
2 detection in genomes and metagenomes**

3 Ashley M. Rooney^{a,b}, Amogelang R. Raphenya^{c,d,e}, Roberto G. Melano^{a,f}, Christine Seah^f, Noelle
4 R. Yee^b, Derek R. MacFadden^g, Andrew G. McArthur^{c,d,e}, Pierre H.H. Schneeberger^{a,b,h} & Bryan
5 Coburn^{a,b,i*}

6 **Affiliations:** ^aDepartment of Laboratory Medicine and Pathobiology, Faculty of Medicine,
7 University of Toronto, Toronto, ON, Canada; ^bUniversity Health Network, Division of Infectious
8 Diseases and Toronto General Hospital Research Institute, Toronto, ON, Canada; ^cDavid Braley
9 Centre for Antibiotic Discovery, McMaster University, Hamilton, ON, Canada; ^dMichael G.
10 DeGroote Institute for Infectious Disease Research, McMaster University, Hamilton, ON,
11 Canada; ^eDepartment of Biochemistry and Biomedical Science, McMaster University, Hamilton,
12 ON, Canada; ^fPublic Health Ontario Laboratory, Toronto, ON, Canada; ^gOttawa Hospital
13 Research Institute, Ottawa, ON, Canada; ^hDepartment of Medical Parasitology and Infection
14 Biology, Swiss Tropical and Public Health Institute, University of Basel, Basel, Switzerland;
15 ⁱDepartment of Medicine, Faculty of Medicine, University of Toronto, Toronto, ON, Canada

16 **Address correspondence to:** Bryan Coburn, University Health Network, Division of Advanced
17 Diagnostics, Princess Margaret Cancer Research Tower, 101 College St., Toronto, ON, M5G 1L7,
18 bryan.coburn@uhn.ca

19 **Alternate corresponding author:** Pierre Schneeberger, Swiss Tropical and Public Health
20 Institute, Department Medical Parasitology and Infection Biology, University of Basel,

21 Socinstrasse 57. 4051 Basel, Petersplatz 1. 4001 Basel, Switzerland,

22 pierre.schneeberger@swisstph.ch

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41

42

43

44 **Abstract**

45 Short-read sequencing provides a culture-independent method for the detection of
46 antimicrobial resistance (AMR) genes from single bacterial genomes and metagenomic samples.
47 However, the performance characteristics of these approaches have not been systematically
48 characterized. We compared assembly- and read-based approaches to determine sensitivity,
49 positive predictive value, and sequencing limits of detection required for AMR gene detection
50 in an *Escherichia coli* ST38 isolate spiked into a synthetic microbial community at varying
51 abundances. Using an assembly-based method the limit of detection was 15X genome
52 coverage. We are confident in AMR gene detection at target relative abundances of 100% to
53 1%, where a target abundance of 1% would require assembly of approximately 30 million reads
54 to achieve 15X target coverage. Recent studies assessing AMR gene content in metagenomic
55 samples may be inadequately sequenced to achieve high sensitivity. Our study informs future
56 sequencing projects and analytical strategies for genomic and metagenomic AMR gene
57 detection.

58

59

60

61

62

63

64

65 **Introduction**

66 Increasing throughput and decreasing costs of DNA sequencing have made whole
67 genome and metagenomic sequencing accessible for antimicrobial resistance (AMR) gene
68 detection on a broad scale. This technology is a useful epidemiological tool^{1,2} and there are
69 increased efforts to correlate isolate genotype with phenotypic resistance^{3,4}. The 'resistome'⁵ is
70 the total genetic content of the microbiome with the potential to confer resistance to
71 antibiotics, and there has been significant interest in characterizing the AMR gene content in
72 the environment^{6–8}, humans^{9,10}, and other mammals^{11,12}. Large trials investigating antibiotic
73 efficacy have also included the development of AMR in the gut microbiome as an outcome¹³.

74 Novel methods for AMR gene detection have been developed to tackle the challenge of
75 AMR gene identification from single isolates and metagenomic samples, generally using either
76 assembly-based or read-based approaches, but there is currently no universal standard¹⁴.
77 Notably, there are no recommendations for optimal sequencing depths required to identify
78 AMR genes in complex metagenomic samples, and the performance characteristics of different
79 sequencing depths using common AMR detection tools for genomes and metagenomes have
80 not been established.

81 In this study, we used the Resistance Gene Identifier (RGI) and the Comprehensive
82 Antibiotic Resistance Gene Database (CARD) and compared an assembly- and read-based
83 approach to assess the limits of detection, sensitivity, and PPV of sequencing to detect known
84 AMR genes in a multidrug- resistant *E. coli* isolate that represented varying abundances in a
85 complex metagenome. We highlight the importance of maintaining minimum target genome
86 coverage to detect AMR genes when the target organism is at varying relative abundances in a

87 metagenomic sample and provide an estimate of minimum required sequencing depths of
88 target organisms to maintain adequate sensitivity.

89 **Results**

90 **Assembly-based AMR gene detection in *Escherichia coli* ST38**

91 To ensure that sequencing effort was not a limiting factor in AMR gene detection, we
92 subjected an *Escherichia coli* (*E. coli*) ST38 isolate to deep sequencing and obtained
93 approximately 136 million reads (~6,800X genome coverage). To simulate sequencing at lower
94 depths, we randomly subsampled 5,000,000 (~250X), 1,000,000 (~50X), 500,000 (~25X),
95 300,000 (~15X), 250,000 (~12.5X), 200,000 (~10X), 150,000 (~7.5X), 100,000 (~5X), 50,000 (~
96 ~2.5X), and 10,000 (~0.5X) read pairs, and bootstrapped each subsample 100 times, with
97 replacement, to provide confidence in the AMR genes detected in each subsample. We
98 considered AMR genes detected with $\geq 90\%$ detection frequency as high confidence genes,
99 whereas those detected with $\leq 50\%$ detection frequency were considered low confidence
100 genes.

101 Using the SPAdes genome assembler, a sequencing depth of 300,000 reads or
102 approximately 15X coverage was sufficient to detect *bla*_{CTX-M-15}, and *parC* and *gyrA* single
103 nucleotide polymorphisms (SNPs) as well as 69 other genes with greater than $\geq 90\%$ detection
104 frequency (Fig. 1a). Other resistance genes included 3 different beta-lactamases (*bla*_{TEM-1},
105 *bla*_{OXA-1}, and *bla*_{AmpC}), 5 unique aminoglycoside transferases, and 46 distinct efflux pump genes
106 (Fig. 1b). A lower sequencing depth was adequate to detect the SNPs with $\geq 90\%$ detection
107 frequency (150,000 reads for *gyrA* and 100,000 for *parC*) compared to *bla*_{CTX-M-15} (200,000
108 reads) (Fig. 1a). There were AMR genes detected less frequently ($\leq 50\%$ detection frequency)

109 across all sequencing depths except at 500,000 reads where no AMR genes were detected at
110 $\leq 50\%$ detection frequency (Fig. 1a).

111 To demonstrate how sequencing depth affects the performance of AMR gene detection,
112 we used a single 5 million read subsample (~250X coverage) as a reference to calculate
113 sensitivity and positive predictive value (PPV) across subsamples. We did not use specificity as a
114 metric to assess performance due to the high number of true negatives which would inflate
115 specificity. A depth of 300,000 reads performed similarly to 1 million reads for sensitivity (1.00
116 ± 0.00 vs 1.00 ± 0.00 , Fig. 1c) and PPV (mean = 1.00 ± 0.00 vs 1.00 ± 0.00 , Fig. 1d) with low false
117 negatives (0.09 ± 0.29 , Fig. 1e) and false positives (0.02 ± 0.14 , Fig. 1f) (mean and standard
118 deviation).

119 The Basic Local Alignment Tool (BLAST) is a highly sensitive alignment tool in common
120 use¹⁵. We aimed to compare high and low confidence AMR genes predicted using BLAST or
121 DIAMOND, a faster alternative sequence alignment tool to BLAST¹⁵. Overall, BLAST predicted
122 more AMR genes across all subsamples (Supplementary Fig. 1c). Using BLAST, as sequencing
123 depth increased, 72 AMR genes achieved $\geq 90\%$ detection frequency by 300,000 reads
124 (Supplementary Fig. 1a), which is consistent with results using DIAMOND (Supplementary Fig.
125 1b). Between BLAST and DIAMOND, the genes predicted with $\geq 90\%$ detection frequency at
126 subsamples $\geq 300,000$ reads were similar in number (approximately 72 genes were detected by
127 both methods) as well as annotation (Supplementary Fig. 1d & 1f). Across all subsamples, more
128 genes were predicted with $\leq 50\%$ detection frequency using BLAST. For example, in the 300,000
129 read subsample, 2 genes were detected with $\leq 50\%$ detection frequency using DIAMOND and 16
130 genes were predicted using BLAST (Supplementary Fig. 1e). Of the total AMR genes detected by

131 BLAST and DIAMOND with $\leq 50\%$ detection frequency, at subsamples 300,000, 500,000,
132 1,000,000 and 5,000,000 reads, 14/16 (87%), 4/4 (100.0%), 8/11 (73%), and 22/25 (88%),
133 respectively, were unique to BLAST (Supplementary Fig. 1f). We proceeded to use DIAMOND
134 for the remainder of the study as fewer low confidence genes were predicted.

135 **Read-based AMR gene detection in *E. coli* ST38**

136 We next compared AMR genes predicted using a read-based approach to AMR genes
137 predicted with an assembly-based approach in the *E. coli* ST38 isolate. We chose to use KMA
138 over other read alignment tools as it produces a consensus sequence which can be used for SNP
139 detection¹⁶. For these analyses, we compared the AMR genes predicted across subsamples
140 using KMA to those predicted using SPAdes assemblies in the 5 million read subsample.

141 We found that 200,000 reads (~10X coverage) were sufficient to identify most AMR
142 genes using KMA read alignment with $\geq 90\%$ detection frequency (Fig. 2a) with a sensitivity of
143 $93\% \pm 0.0\%$ and 5.0 ± 0.2 false negatives (mean \pm standard deviation) (Fig 2c & 2e). However,
144 there were genes that had $\leq 50\%$ detection frequency at 200,000 reads including the known
145 *gyrA* and *parC* SNPs as well as *KpnE*, *KpnF*, *bla_{TEM-1}*, and the gene annotated as homologous to
146 *Haemophilus* penicillin-binding protein 3 conferring resistance to beta-lactams. The AMR gene
147 *rsmA* required 5 million reads to achieve a detection frequency of $\geq 90\%$ (Fig. 2a). At 200,000
148 reads there were 30 additional AMR genes predicted at $\geq 90\%$ detection frequency that were
149 not detected in the reference (Fig. 2b), with a PPV of $66\% \pm 0.0\%$ and 34.8 ± 1.1 false positives
150 at this depth (Fig. 2d & 2f).

151 When a read-based approach is used to identify AMR genes, strategies are often applied
152 to improve precision and remove false positives. However, there is a trade-off between

153 increased precision and decreased recall which is important to quantify. We assessed the
154 effects of four AMR gene filtering strategies, at a range of cut-off values, on the performance of
155 the 200,000 read depth for the detection of AMR genes in *E. coli* ST38 using KMA. The AMR
156 gene filtering strategies included percent coverage, average depth of coverage, number of
157 completely mapped reads, and the average mapping quality score (MAPQ score). The unfiltered
158 precision and recall was 66% and 93%, respectively. No filtering strategy that we assessed
159 significantly improved the precision, based on the precision-recall curve (precision = PPV; recall
160 = sensitivity, Fig. 3). Out of the four strategies, percent coverage achieved the greatest increase
161 in precision at the highest stringency cut-off of 100% allele coverage (unfiltered precision: 66%;
162 filtered precision: 76%). However, at this cut-off, the recall decreased from 93% to 79%. When
163 filtering by depth of coverage, completely mapped reads, or MAPQ score, the highest increase
164 in precision was from 66% to 74% (recall: 76%), 74% (recall: 54%), and 69% (recall: 87%),
165 respectively. As expected, as the cut-offs became more stringent for each filtering strategy, the
166 recall decreased. Filtering based on percent coverage had the least affect on recall, even at the
167 highest stringency cut-off (100% allele coverage) compared to the other strategies.

168 **Detection of *E. coli* ST38 AMR genes at a range of relative abundances in a complex
169 metagenomic sample**

170 To demonstrate the effect of target organism relative abundance on AMR gene
171 detection in a multi-species metagenome consisting of 34 bacterial species, the DNA of *E. coli*
172 ST38 and the complex community were combined to create synthetic metagenomic samples
173 where *E. coli* ST38 represented approximately 90%, 50%, 10%, and 1% of the total
174 metagenome. Based on the sequencing limit of detection of 300,000 reads in the single *E. coli*

175 ST38 isolate (100% relative abundance), we estimated that at 90%, 50%, 10%, and 1% relative
176 abundance, 333,333, 600,000, 3,000,000 and 30,000,000 reads, respectively, would be required
177 to detect the known AMR genes (*bla*_{CTX-M-15}, *parC*, *gyrA* SNPs) contributed by the *E. coli* ST38
178 isolate with \geq 90% detection frequency.

179 Using metaSPAdes as the metagenomic assembly tool, we observed that as the *E. coli*
180 ST38 relative abundance decreased, the number of reads necessary to detect the *bla*_{CTX-M-15}
181 (Fig. 4a), the *gyrA* SNPs (Fig. 4b), and the *parC* SNP (Fig. 4c), as well as the 5 aminoglycoside
182 transferases, *bla*_{TEM-1}, and *bla*_{OXA-1}, increased (Supplementary Figure 2a-f). The detection rate
183 approximated our expectations at relative abundances $>1\%$ (Fig. 4a-c). For the combined
184 sample containing *E. coli* ST38 at 1% relative abundance, we did not detect the *gyrA* SNPs in any
185 of the 10 bootstraps at 30,000,000 reads (Fig. 4b), while the *bla*_{CTX-M-15} and the *parC* SNP had a
186 detection frequency of 90% (9/10 bootstraps) and 100% (10/10 bootstraps), respectively (Fig.
187 4a & 4c).

188 Although read-based approaches are often thought of as highly sensitive tools for AMR
189 gene detection in low abundance organisms¹⁴, KMA did not significantly improve the limit of
190 detection for *bla*_{CTX-M-15} (Fig. 4d), the 5 aminoglycoside transferases or *bla*_{OXA-1} (Supplementary
191 Figure 3a-e & 3g) compared to assembly. The detection frequency was low for the *gyrA* and
192 *parC* SNPs (Fig. 4e & 4f), as well as *bla*_{TEM-1} (Supplementary Figure 3f) across subsamples and
193 metagenomic samples. Instead of *bla*_{TEM-1}, KMA aligned reads to three TEM-variants with \geq 90%
194 detection frequency in at least one metagenomic sample, of which *bla*_{TEM-181} had the most
195 reads aligning to the allele (Supplementary Figure 4d-f). In comparison, all TEM-variants
196 predicted using KMA had low detection frequency using metaSPAdes (Supplementary Figure

197 4a-c). Sanger sequencing confirmed the presence of *bla_{TEM-1}* and not *bla_{TEM-181}* in the *E. coli* ST38
198 isolate. Lastly, the total number of high confidence AMR genes in the metagenomic samples
199 was significantly higher using KMA compared to those genes detected using metaSPAdes.

200 **Validation of 15X coverage across *E. coli* isolates**

201 To validate the 300,000 read depth/15X coverage threshold, we applied our assembly-
202 based approach to 948 *E. coli* isolates¹⁷. The isolates were previously sequenced to an average
203 of 100X coverage using 150-bp paired-end Illumina sequencing. We performed a subsample
204 from each isolate at 300,000 reads and compared the AMR genes predicted at 300,000 reads to
205 the AMR genes predicted at the original sequencing depth to calculate sensitivity, PPV, and F1
206 score for each isolate. The F1 score is a harmonic measure of sensitivity and PPV, where a score
207 of 1 would indicate perfect sensitivity and PPV.

208 Across the *E. coli* isolate set, we observed a total of 322 unique AMR genes.
209 Performance of the 300,000 read depth is summarized in Figure 5. The F1 score was 1 for
210 658/948 (69.4%) isolates. There were 290/948 (30.6%) isolates with a F1 score of <1, where
211 228/290 (78.6%) isolates had an F1 score between 0.99 – 0.98, 49/290 (16.9%) had an F1 score
212 between 0.95 – 0.97, 11/290 (3.8%) had an F1 score between 0.90 – 0.94 and the remaining
213 isolates had an F1 score 0.89 and 0.65. Of the 290 isolates with F1 score <1 (less than perfect
214 agreement), 84 (29.0%) had a PPV of <1 and 261 (90.0%) had a sensitivity of <1. For the isolates
215 with a PPV of <1, the median number of false positives was 1 (range 1-70). The isolate with 70
216 false positives is an outlier. For the isolates with a sensitivity of <1, the median number of false
217 negatives was 1 (range 1-15).

218 Of the 322 unique AMR genes, 21 genes (6.5%) were classified as true positive for all
219 948 isolates, where 90.5% (19/21) were efflux-associated genes. The top three AMR genes that
220 contributed the most false negatives were *APH(6)-Id* (n = 25), *sul2* (n = 23), and *mphA* (n = 23),
221 while the top three genes that contributed the most false positives were *bla_{OXA-320}* (n = 13),
222 *aadA* (n = 6), and *bla_{OXA-140}* (n = 6).

223 **Validation of 15X coverage in metagenomic samples**

224 From a public dataset of 10 rectal surveillance swabs which were vancomycin resistant
225 *Enterococcus* (VRE) positive by culture and *vanA* positive in 9/10 swabs by Illumina sequencing
226 ¹⁸, we validated 15X *Enterococcus* genome coverage for the detection of *vanA*. The study
227 authors performed 2 X 75 bp sequencing and achieved a mean 9.1 million reads (range: 5.7
228 million – 15 million reads), post-quality filtering and removal of human reads. The rectal swab
229 samples had a range of *Enterococcus* relative abundances (median: 0.10; range: 0.80 – 0.0002)
230 and genome coverages (median: 21X; range: 375X – 0.07X) (Fig. 6a). Rectal swab number 8 had
231 the highest *Enterococcus* relative abundance of 0.80 and, due to the large number of
232 sequencing reads (125 million), had the largest estimated target genome coverage of ~375X.
233 Rectal swab number 4 had the lowest *Enterococcus* relative abundance (0.0002) and 91 million
234 sequences, which resulted in a target genome coverage of ~0.07X for this sample.

235 To assess whether a minimum of 15X target genome coverage is sufficient to detect
236 *vanA* in the rectal swab metagenomes, we subsampled reads to achieve a range of target
237 genome coverages from 0.5X - 15X and then bootstrapped 10 times at each subsample to
238 determine *vanA* detection frequency. The results of this analysis are shown in Figure 6b. To
239 achieve 100% detection frequency of the *vanA* gene across rectal swab samples, 5 samples

240 required *Enterococcus* genome coverage of less than 5X (rectal swabs 1,5-7, and 10), while 2
241 required at least 15X coverage (rectal swabs 3 and 8). At 15X *Enterococcus* genome coverage,
242 *vanA* was detected in 10/10 bootstraps for all samples that had adequate sequencing depth for
243 subsampling. Rectal swab number 4 did not have enough reads to achieve 0.5X *Enterococcus*
244 genome coverage and *vanA* was not detected when we analyzed all reads available, which is
245 consistent with the authors' published findings that describe their inability to detect *vanA* using
246 paired-end Illumina sequencing¹⁸.

247 **Estimates of the sensitivity of sequencing depth for AMR gene detection in published data**

248 **sets**

249 Recent publications assessing AMR gene content in metagenomic samples may not have
250 achieved optimal sensitivity for AMR gene detection if they were to use a contig assembly
251 approach. As we have observed, the relative abundance of the target organism affects
252 sensitivity to detect AMR genes in a metagenomic sample. We gathered sequencing depths and
253 read length information from three recently published studies that reported AMR genes in
254 metagenomic samples. Study 1 assessed and compared the resistome of 1174 gut and oral
255 samples from previously published sources distributed by country⁹. For our analyses, we
256 included 1132/1174 from Study 1 for which we had complete read length data (excluding
257 42/1174, 3.6%). Study 2 performed a longitudinal assessment of the gut microbiota and
258 resistome of healthy veterinary students exposed to a Chinese swine farm environment. A total
259 63 metagenomic samples were sequenced which consisted of human stool and environmental
260 samples⁶. Study 3 was conducted in Denmark and evaluated the changes in the gut microbiota
261 composition and resistome of 12 healthy male volunteers before and after antimicrobial

262 exposure¹⁰. A total of 57 stool samples were subject to metagenomic analyses. Studies 1 and 3
263 used a read-based approach for AMR gene prediction, while Study 2 used an assembly-based
264 approach.

265 We were interested in estimating AMR gene detection frequency from the published
266 sample sequencing depths for a hypothetical target organism (presumed genome size 6 Mbp)
267 at a range of potential relative abundances. Assuming detection frequency is related to
268 sequencing sensitivity, we calculated coverage as an estimate of sequencing depth and
269 interpolated detection frequency values from a sigmoidal curve fit to the *E. coli* ST38 *bla*_{CTX-M-15}
270 detection data as seen in Figure 4a.

271 As the relative abundance of the hypothetical target organism decreased, more
272 sequencing effort was required to achieve 100% estimated detection frequency of all AMR
273 genes (Fig. 7). Most published samples had achieved ≥95% estimated detection frequency for
274 all AMR genes for a target organism at relative abundance of 100% (1251/1252; 99.9%), 90%
275 (1250/1252; 99.8%) and 50% (1247/1252; 99.6%). However, the proportion of samples with at
276 least 90% estimated detection frequency was lower for a target organism relative abundance of
277 10% (1090/1252; 87.1%) and 1% (454/1252; 36.3%). Additionally, 29.5% (369/1252) of samples
278 were not sequenced sufficiently to achieve >50% estimated detection frequency for a target
279 organism relative abundance of 1%, where 9.2% (115/1252) had less than 1% estimated
280 detection frequency (Fig. 7), suggesting that in these studies, sequencing depth may be
281 inadequate to achieve a high sensitivity for detection of AMR genes in low abundance
282 organisms.

283 **Discussion**

284 Our goal was to characterize the performance (sensitivity, PPV, and limits of detection)
285 of genomic and metagenomic approaches for the detection of known predictors of AMR in
286 microbes to inform the use of these approaches in human, animal, and environmental studies.
287 It is axiomatic that sequencing depth affects AMR gene assay sensitivity in single isolates^{19,20}
288 and within a microbiome²¹, however, the performance characteristics of sequencing have not
289 been systematically assessed. In published reports, a range of whole genome sequencing
290 depths for single isolates, from 30X coverage up to 100X coverage, are often used to define
291 quality control limits, but these are not considered standard^{3,19,22,23}. Estimating the coverage of
292 the metagenome required to ensure high sensitivity is not a new concept²⁴, but we are
293 unaware of a study that attempts to precisely quantify sequencing depths required to detect
294 AMR genes in a target organism across varying relative abundances in mixed metagenomes
295 using standard methods.

296 Using a *de novo* contig-assembly approach, we found that approximately 15X coverage
297 (300,000, 2 X 150 bp paired-end reads of a 6 Mbp genome) provides similar sensitivity to higher
298 sequencing depths for the detection of AMR genes in *E. coli* isolates and is sufficient for
299 detecting SNPs and other resistance genes. Although sequencing depths as low as 0.5 million
300 reads have been proposed to capture the total compositional information of metagenomes²⁵,
301 greater sequencing depth is required for the detection of AMR genes in organisms with low
302 relative abundance, which can require as many as 30 million reads to achieve adequate
303 sensitivity for organisms at a relative abundance of 1%.

304 For some study purposes, detection of AMR genes in low abundance organisms may be
305 critical for study interpretation. Human observational studies have demonstrated that

306 pathogens at both high and low relative abundances in complex gut microbial communities are
307 associated with subsequent infections or death. Dominance of a microbial community by a
308 pathogen is associated with subsequent infection^{26–28}, but even at relative abundances as low
309 as 1% - 0.1%, pathogens detected in stool have been implicated in subsequent bacteremia in
310 hematopoietic stem cell transplant recipients²⁹, as well as bacteriuria and urinary tract
311 infection³⁰, indicating that detection of AMR genes may be clinically significant even at very low
312 relative abundance thresholds. Based on our findings, approximately 64% of the samples in
313 recent studies evaluating AMR gene content in the metagenome are not sequenced at a
314 sufficient depth to detect AMR genes in a target organism at 1% relative abundance. Thus,
315 potentially clinically meaningful resistance determinants may not be detected with common
316 sequencing depths such as those we analyzed in published studies.

317 Assembly is time-consuming, requires large amounts of computing power for
318 metagenomic samples, and may also contribute to loss of data³¹. Alternative approaches to
319 assembly such as read alignment^{16,32–34} and kmer-based approaches³¹ may require less
320 sequencing information for AMR gene detection, which is useful for detecting AMR genes in
321 low abundance organisms in complex communities. Compared to assembly, our read-based
322 approach (KMA) did not significantly improve the limit of detection in *E. coli* ST38 or
323 metagenomes, even for a low abundance target, and at times suffered from the AMR allele
324 network problem³⁵, where reads from a single gene were aligned to multiple closely related
325 reference alleles, e.g. *bla_{TEM-181}*, despite its improvements over other read alignment tools for
326 this issue¹⁶. Additionally, some genes (e.g. *gyrA* and *parC* SNPs and *bla_{TEM-1}*) that had a high
327 detection frequency with assembly, had a low detection frequency with KMA and there were a

328 large number of false positive genes that could not be filtered without affecting the overall
329 sensitivity of the assay. Yet, our comparison of BLAST and DIAMOND with open reading frames
330 predicted from assembled contigs illustrated that attention should also be paid to local
331 alignment algorithm choice when using assembly-based approaches, as DIAMOND generated
332 less low confidence predictions.

333 A large majority of the AMR genes detected in the *E. coli* ST38 isolate and across the *E.*
334 *coli* isolate set were efflux-associated AMR genes. Efflux-associated genes are of uncertain
335 relevance for prediction of microbial phenotypes^{36,37}. In some scenarios (e.g. human infections)
336 it may be advantageous to limit AMR gene detection to acquired mechanisms of resistance by
337 using a database such as ResFinder³⁸, or by filtering efflux-associated genes from the total AMR
338 genes detected, which has been previously shown to improve predictions of isolate
339 antimicrobial susceptibility using CARD³⁷. Accurate prediction of plasmids, often associated
340 with clinically important AMR genes, remains difficult as short-read Illumina sequencing
341 provides highly accurate base calling, but repetitive DNA regions complicate genome assembly
342 resulting in fragmented, short contigs³⁹. AMR genes may be split between multiple contigs⁴⁰
343 leaving plasmid sequences obscured⁴¹. Long-read sequencing technology provides a promising
344 alternative to short-read sequencing and can overcome the issue of fragmented contig
345 assembly³⁹.

346 Our approach has the following limitations. We modelled a single approach utilizing a
347 widely used sequencing strategy, two bioinformatic pipelines and one AMR detection platform
348 (CARD) for a single organism (*E. coli*). These selections were made to reflect dominant modes of
349 metagenome analysis in a clinically relevant organism to define the 'order of magnitude' of

350 depth required for AMR gene detection from metagenomes, which may not be generalizable to
351 all organisms, community types or modes of resistance. A main limitation of AMR gene
352 prediction from sequencing data is the chosen database that can potentially increase false
353 negatives. However, CARD is widely used, updated on a monthly basis, and is representative of
354 known AMR gene diversity, especially for well-characterized pathogens such as *E. coli*⁴². Human
355 metagenomic samples often have human DNA that can account for a large proportion of the
356 total sample, which impacts sequencing strategies^{43,44}. An understanding of the total genetic
357 material contributed by human reads prior to sequencing would further inform sequencing
358 effort required to maintain a minimum sequencing depth for AMR gene detection.

359 We have quantified sequencing depths needed to detect AMR genes in *E. coli* whole
360 genomes and in *E. coli* from high to low relative abundances among a complex community. A
361 minimum of 15X coverage is needed for the detection of AMR genes in *E. coli* using our AMR
362 gene identification approach. For metagenomic samples, 15X coverage is also sufficient to
363 detect known AMR genes in *E. coli*, but the number of sequences must increase proportionally
364 to the decrease in relative abundance of the target organism. We believe that this analysis
365 provides a robust benchmarking of sequencing effort for metagenomic studies in which
366 detection of resistance is a specified outcome.

367 **Methods**

368 **Sample preparation and sequencing**

369 From a collection of previously characterized *E. coli* isolates¹⁷, we selected a multidrug-
370 resistant *E. coli* ST38 with an extended spectrum beta-lactamase (*bla*_{CTX-M-15}) and resistance-
371 conferring SNPs in *parC* (S80I) and *gyrA* (S83L, D87N). Briefly, *E. coli* ST38 isolate was cultured

372 from a glycerol stock on LB agar and a single colony was inoculated into 25 ml of LB broth,
373 which was placed on a shaker incubator (130 rpm) at 37°C for 4 hours until media was turbid.
374 Turbid media (25 ml) was transferred to a 50 ml conical tube, subject to centrifugation at 2500
375 g, the supernatant removed, and the pellet re-suspended in 500 µl of LB broth. A description of
376 the complex community (MET-1) preparation was described previously⁴⁵. Aliquots of MET-1
377 were stored at -80°C prior to use.

378 DNA was extracted from thawed MET-1 (250 µl) and the *E. coli* isolate in LB broth (250
379 µl) using the DNeasy PowerSoil kit (Qiagen) and DNA concentration was measured using a
380 Qubit Fluorometer (Thermo Fisher), following the manufacturer's instructions, respectively. *E.*
381 *coli* and MET-1 DNA were combined to a final concentration of 20.1 ng/µl, while varying the
382 concentration of *E. coli* so that it approximately represented 90%, 50%, 10%, 1%, 0.1%, 0.01%,
383 0.001%, and 0.0001% relative to MET-1. Sequencing libraries were prepared using the Nextera
384 DNA Flex kit (Illumina) following the manufacturer's instructions and stored at -20°C. All 10
385 samples (the *E. coli* ST38 isolate, MET-1, and 8 combined samples) were subject to paired-end
386 sequencing at 2 X 150 bp on the NovaSeq 6000 at the Princess Margaret Genomics Centre.
387 Since we did not achieve the minimum sequencing depth needed to detect AMR genes in the
388 combined samples where *E. coli* ST38 represented 0.1%, 0.01%, 0.0001%, and 0.0001% we did
389 not analyze these samples in our study.

390 **Bioinformatic analyses**

391 From each pair of fastq files, Seqtk⁴⁶ was used to subsample *n* number of reads. At each
392 subsample, bootstrapping was performed (sampling with replacement) 100 times unless
393 otherwise stated, where all 100 bootstraps of a subsample had a unique seed number to ensure

394 every bootstrap was a random sampling of reads. Paired-end fastq files (read 1 and read 2)
395 were assessed for quality using FastQC⁴⁷. Nextera adapters were removed with Trimmomatic⁴⁸
396 v.0.39. Reads for *E. coli* genomes as well as MET-1 and the combined samples were assembled
397 into contigs using SPAdes⁴⁹ v.3.13.1, specifying the *--careful* flag, and metaSPAdes⁵⁰ v.3.13.1,
398 respectively, using the recommended kmer lengths 21, 33, 55, and 77. We used Metaphlan2⁵¹
399 v.2.9.21 to confirm sample taxonomy, including the identity of all *E. coli* isolates and the
400 relative abundance of *Enterococcus* species in the validation sets, respectively.

401 To predict AMR genes from contigs, we used RGI *main* v.5.1.0 of the Comprehensive
402 Antibiotic Resistance Database on default settings (perfect and strict hits identified only)⁴². We
403 specified DIAMOND¹⁵ v.0.8.36, or BLAST⁵² v.2.9.0 (where stated) to perform local alignment of
404 Prodigal-predicted genes within contigs against CARD v.3.1.0^{42,53}. For metagenome assembled
405 contigs, we specified the *--low_quality* flag in RGI *main* to allow prediction of partial open
406 reading frames by Prodigal.

407 To predict AMR genes from raw reads, we used KMA¹⁶ v.1.3.8 within RGI *bwt* v.5.2.0 to
408 align reads to CARD. To predict AMR-conferring SNPs in the *parC* (S80I) and *gyrA* (S83L, D87N)
409 genes, we extracted the consensus sequences generated from these read alignments and used
410 RGI *main* v.5.2.0.

411 **Quantification of AMR genes**

412 To quantify the occurrence of the *parC* (S80I) and *gyrA* (S83L, D87N) SNPs in
413 bootstrapped subsamples of the *E. coli* ST38 isolate and combined samples, we extracted the
414 individual accession numbers (from contig results only) and SNP information (contigs and raw
415 reads) for each gene from all RGI or KMA output files. The *parC* and *gyrA* SNPs were considered

416 present if the mutated amino acid residues S80I in *parC* and S83L, D87N in *gyrA* were correctly
417 predicted by RGI. If RGI predicted other mutated amino acid residues in the *parC* and *gyrA*
418 genes or did not identify any mutated amino acid residues in the *parC* and *gyrA* gene, we did
419 not consider these SNPs as present. For all other resistance genes, we extracted unique genes
420 from the “Best_Hit_ARO” (contigs) or “ARO_Term” (raw reads) column of each sample RGI
421 output file, to create a new “unique AMR genes” file for each sample. Then, using Metaphlan2
422 v.2.9.14 we used the `merge_metaphlan_tables.py` command to merge the “unique AMR genes”
423 files together, where the first column outlined the AMR genes predicted for all samples and the
424 first row indicated the sample names. AMR gene presence was indicated by 1 and absence
425 indicated by 0. Merging the RGI output files allowed us to quantify the frequency at which
426 individual AMR genes were present across samples. We considered an AMR gene present in a
427 bootstrap if the gene occurred at least once. The number of AMR genes present across
428 bootstrapped subsamples was visualized with rarefaction curves and plotted using GraphPad
429 Prism version 9.

430 **Coverage estimation**

431 Sequencing coverage was estimated using the Lander-Waterman equation⁵⁴. For *E. coli*
432 ST38, we assumed a genome size of 6 Mbp. To estimate the number of reads required to detect
433 *E. coli* ST38 at a range of relative abundances, the minimum read requirement (300,000 reads)
434 was divided by the target relative abundance. For example, if the target relative abundance was
435 10%, 300,000/0.10 would equal 3,000,000 reads.

436 **Performance analyses**

437 The performance of AMR gene classification was calculated using sensitivity and positive
438 predictive value. For the single *E. coli* ST38 isolate, the AMR genes predicted in each
439 bootstrapped subsample using a contig assembly approach with SPAdes or a read-based
440 approach with KMA were compared to the AMR genes predicted in a subsample assembled into
441 contigs at 5 million reads (reference). If the AMR gene was present in the bootstrap and
442 reference, this gene was considered a true positive. If an AMR gene was not present in neither
443 the bootstrap nor the reference, this gene was considered a true negative. False positive AMR
444 genes were present in the bootstrap but absent in the reference and false negative AMR genes
445 were absent in the bootstrap but present in the reference. For each bootstrap sample, the true
446 positives, true negatives, false positives, and false negatives were summed and sensitivity and
447 PPV were calculated.

448 **Validation from external datasets**

449 To validate the performance of a 300,000 read depth across a set of *E. coli* isolates¹⁷, we
450 subsampled 300,000 reads, once, from each isolate, assessed quality with FastQC and discarded
451 isolates that failed per base sequence quality. We then compared the AMR genes detected at
452 300,000 reads to the AMR genes detected from the original sequence depth. We summed the
453 true positives, true negatives, false positives, and false negatives for each isolate, then
454 calculated sensitivity, PPV and F1 score as a balanced measure of sensitivity and PPV.

455 To validate 15X target genome coverage in metagenomic samples¹⁸ and to demonstrate
456 *vanA* detection frequency across a range of *Enterococcus* genome coverages (0.5X – 15X), we
457 subsampled each metagenomic sample and bootstrapped each subsample 10 times. Each
458 subsample depth was calculated using the Lander-Waterman equation, as described above,

459 while accounting for the *Enterococcus* relative abundance in the sample, as determined using
460 Metaphlan2. We assumed an *Enterococcus* genome size of 4 Mbp.

461 **AMR gene detection frequency assessment of published datasets**

462 We extracted the sequence depths after quality processing that were provided in each
463 studies' supplementary material for Study 1⁹, 2⁶ and 3¹⁰. For Study 3, we used the sequences
464 reported under the heading "After human contamination removal" under the sub-heading
465 "read-pairs". Sequencing read lengths were reported in Study 2 (2 X 150 base pairs) and 3 (2 X
466 100 base pairs), but for Study 1 we extracted the read lengths from the individual studies
467 referenced within the paper. We then estimated coverage for each sample from the published
468 datasets and for each subsample performed on the samples where *E. coli* ST38 represented
469 100%, 90%, 50%, 10%, and 1% relative abundance, assuming a genome length of 6 Mbp, then
470 log-transformed these values. Using GraphPad Prism version 9.1.2, sigmoidal curves were fit to
471 detection frequency data for the *bla*_{CTX-M-15} for each sample where *E. coli* ST38 represented
472 100%, 90%, 50%, 10%, and 1% relative abundance. The equations were constrained at 0 and
473 100 and detection frequency was interpolated for relative abundances 100%, 90%, 50%, 10%,
474 and 1% based on coverage estimation.

475 **Data availability**

476 The dataset generated during the current study are available in the NCBI sequence read archive
477 under the accession number PRJNA649958.

478 **Competing interests**

479 The authors declare that they have no competing interests.

480 **Author Contributions**

481 AMR, PHHS, and BC conceived and designed the study. BC supervised the overall study. AMR
482 with the guidance and supervision of PHHS performed sample preparation and bioinformatic
483 analyses. ARR provided feedback on and performed some bioinformatic analyses. RM and CS
484 provided qPCR and Sanger sequencing support. RM, DM, and AM provided feedback during
485 analyses. NY provided help with the bioinformatic analyses. AMR generated the figures and
486 wrote the manuscript. All authors provided feedback during manuscript preparation and have
487 read and approved the final manuscript.

488 **Acknowledgments**

489 We greatly appreciate Dr. Emma Allen-Vercoe from the University of Guelph for sending us
490 aliquots of the complex community, MET-1. This research was funded by the Canadian
491 Institutes of Health Research (PJT-156214 to AM) and a David Braley Chair in Bioinformatics to
492 AM. Some computer resources were supplied by the McMaster Service Lab and Repository
493 computing cluster, funded in part by grants to AM from the Canada Foundation for Innovation
494 (34531) and hardware donations for Cisco Systems Canada, Inc.

495

496 **Figure Legends**

497 **Figure 1. a**, A rarefaction plot of *Escherichia coli* ST38 AMR genes detected across subsamples
498 using an assembly-based approach. Individual dots represent a single AMR gene and are
499 connected by lines to demonstrate trends in detection across subsamples. *bla_{CTX-M-15}*, *gyrA* SNPs
500 (S83L, D87N), and *parC* SNP (S80I) are highlighted as previously identified resistance
501 determinants for this strain. The horizontal dotted line marks 90% detection frequency. The red
502 vertical dashed line marks the subsample at 300,000 reads. **b**, Histogram of the number of

503 unique AMR genes with $\geq 90\%$ detection frequency summarized by categories detected across
504 subsamples. **(c-f)** Performance of AMR gene classification across subsamples. A 5 million read
505 subsample was used as reference to calculate sensitivity **(c)**, positive predictive value **(d)**, false
506 negatives **(e)**, and false positives **(f)**. **c-f**, the mean and standard deviation are plotted.

507

508 **Figure 2. a-b**, Rarefaction plots of *Escherichia coli* ST38 AMR genes detected across subsamples
509 using a read-based approach. Individual dots represent a single AMR gene and are connected
510 by lines to demonstrate trends in detection across subsamples. *bla_{CTX-M-15}*, *gyrA* SNPs (S83L,
511 D87N), and *parC* SNP (S80I) are highlighted as previously identified resistance determinants for
512 this strain. The horizontal dotted line marks 90% detection frequency. The red vertical dashed
513 line marks the subsample at 200,000 reads. **a**, AMR genes detected across subsamples which
514 are present in the reference, **b**, AMR genes detected across subsamples which are not present
515 in the reference. **(c-f)** Performance of AMR gene classification across subsamples, where
516 performance was measured by sensitivity(**c**), positive predictive value (**d**), false negatives (**e**),
517 and false positives (**f**). **c-f**, the mean and standard deviation are plotted. **a-e**, the reference used
518 was all AMR genes detected using an assembly-based approach in the *E. coli* ST 38 isolate
519 subsampled at 5 million reads.

520

521 **Figure 3.** Precision-recall curve for an assessment of the effects of filtering strategies on the
522 read-based (KMA) classification of all AMR genes in *Escherichia coli* ST38 at 200,000 reads.
523 Precision is a measure of positive predictive value and recall is a measure of sensitivity.

524

525 **Figure 4.** Detection of *bla*_{CTX-M-15} (a,d), *gyrA* (S83L, D87N) (b,e), and *parC* (S80I) (c,f) from *E. coli*
526 ST38 isolate (100%) across subsamples at varying strain relative abundances (90%, 50%, 10%,
527 1%) in a complex (multi-species) metagenomic sample using either an assembly-based
528 approach (metaSPAdes)(a-c) or a read-based approach (KMA)(d-f). The horizontal dotted line
529 marks 90% detection frequency. The vertical dotted lines represent the read depths where each
530 gene was estimated to be detected with $\geq 90\%$ detection frequency for each colour matched
531 relative abundance using an assembly-based approach. The estimated read depths are 300,000,
532 333,333, 600,000, 3,000,000, and 30,000,000 to detect *bla*_{CTX-M-15}, *gyrA*, and *parC* SNPs in $\geq 90\%$
533 detection frequency when *E. coli* ST38 is at 100%, 90%, 50%, 10%, and 1% relative abundance
534 respectively. For all subsamples 100 bootstraps were performed, except for the 30,000,000
535 read subsample where 10 bootstraps were performed, which is marked by an asterisk (*). The
536 30,000,000 read subsample was performed on one combined sample where *E. coli* represented
537 1% relative abundance.

538

539 **Figure 5.** Relative frequency distribution of *E. coli* isolates (n = 948) by performance of 300,000
540 reads for AMR gene detection using an assembly-based approach. Performance is measured by
541 sensitivity, positive predictive value (PPV), and F1 score.

542

543 **Figure 6.** Detection of *vanA* in rectal swab samples positive for vancomycin-resistant
544 *Enterococcus* from a public dataset. **a** *Enterococcus* relative abundance by estimated total
545 genome coverage, each rectal swab sample is represented by an icon. **b** *vanA* detection
546 frequency across genome coverages for each rectal swab sample. Rectal swab sample 4 is not

547 plotted as *vanA* was not detected with the total number of sequences available. For each
548 sample, 10 bootstraps were performed at each genome coverage depth. To calculate
549 *Enterococcus* genome coverage, we assumed a genome size of 4 Mbp.

550

551 **Figure 7.** Estimated AMR gene detection frequency by sample coverage of a hypothetical target
552 organism at different relative abundances. Sample sequencing information was extracted from
553 three published datasets (n = 1252). Detection frequency was estimated by interpolating values
554 from sigmoidal curves fit to the *bla*_{CTX-M-15} data from samples where *E. coli* represented 100%,
555 90%, 50%, 10%, and 1% relative abundance. Coverage was estimated assuming a target
556 genome size of 6 Mbp.

557

558 **References**

- 559 1. Jackson, B. R. *et al.* Implementation of Nationwide Real-time Whole-genome Sequencing
560 to Enhance Listeriosis Outbreak Detection and Investigation. *Clin. Infect. Dis.* **63**, 380–386
561 (2016).
- 562 2. Harrison, L. H. *et al.* Outbreak of Vancomycin-resistant *Enterococcus faecium* in
563 Interventional Radiology: Detection Through Whole-genome Sequencing-based
564 Surveillance. *Clin. Infect. Dis.* **70**, 2336–2343 (2020).
- 565 3. Macfadden, D. R. *et al.* Comparing Patient Risk Factor-, Sequence Type-, and Resistance
566 Locus Identification-Based Approaches for Predicting Antibiotic Resistance in *Escherichia*
567 *coli* Bloodstream Infections. *J. Clin. Microbiol.* **57**, e01780-18 (2019).
- 568 4. Břinda, K. *et al.* Rapid inference of antibiotic resistance and susceptibility by genomic

569 neighbour typing. *Nat. Microbiol.* **5**, 455–464 (2020).

570 5. D’Costa, V. M., McGrann, K. M., Hughes, D. W. & Wright, G. D. Sampling the antibiotic
571 resistome. *Science*. **311**, 374–377 (2006).

572 6. Sun, J. *et al.* Environmental remodeling of human gut microbiota and antibiotic resistome
573 in livestock farms. *Nat. Commun.* **11**, 1427 (2020).

574 7. Chng, K. R. *et al.* Cartography of opportunistic pathogens and antibiotic resistance genes
575 in a tertiary hospital environment. *Nat. Med.* **26**, 941–951 (2020).

576 8. Pehrsson, E. C. *et al.* Interconnected microbiomes and resistomes in low-income human
577 habitats. *Nature* **533**, 212–216 (2016).

578 9. Carr, V. R. *et al.* Abundance and diversity of resistomes differ between healthy human
579 oral cavities and gut. *Nat. Commun.* **11**, 693 (2020).

580 10. Palleja, A. *et al.* Recovery of gut microbiota of healthy adults following antibiotic
581 exposure. *Nat. Microbiol.* **3**, 1255–1265 (2018).

582 11. Lim, S.-K., Kim, D., Moon, D.-C., Cho, Y. & Rho, M. Antibiotic resistomes discovered in the
583 gut microbiomes of Korean swine and cattle. *Gigascience* **9**, 1–11 (2020).

584 12. Liu, J. *et al.* The fecal resistome of dairy cattle is associated with diet during nursing. *Nat.*
585 *Commun.* **10**, 4406 (2019).

586 13. Doan, T. *et al.* Gut microbiome alteration in MORDOR I: a community-randomized trial of
587 mass azithromycin distribution. *Nat. Med.* **25**, 1370–1376 (2019).

588 14. Boolchandani, M., D’Souza, A. W. & Dantas, G. Sequencing-based methods and resources
589 to study antimicrobial resistance. *Nat. Rev. Genet.* **20**, 356–370 (2019).

590 15. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.

591 *Nat. Methods* **12**, 59–60 (2015).

592 16. Clausen, P. T. L. C., Aarestrup, F. M. & Lund, O. Rapid and precise alignment of raw reads
593 against redundant databases with KMA. *BMC Bioinformatics* **19**, 307 (2018).

594 17. Macfadden, D. R. *et al.* Using Genetic Distance from Archived Samples for the Prediction
595 of Antibiotic Resistance in *Escherichia coli*. *Antimicrob. Agents Chemother.* **64**, e02417-19
596 (2020).

597 18. Yee, R. *et al.* Metagenomic next-generation sequencing of rectal swabs for the
598 surveillance of antimicrobial-resistant organisms on the Illumina MiSeq and Oxford
599 MinION platforms. *Eur. J. Clin. Microbiol. Infect. Dis.* **40**, 95–102 (2021).

600 19. Doyle, R. M. *et al.* Discordant bioinformatic predictions of antimicrobial resistance from
601 whole-genome sequencing data of bacterial isolates: an inter-laboratory study. *Microb.*
602 *Genomics* **6**, e000335 (2020).

603 20. Cooper, A. L. *et al.* Systematic Evaluation of Whole Genome Sequence-Based Predictions
604 of *Salmonella* Serotype and Antimicrobial Resistance. *Front. Microbiol.* **11**, 549 (2020).

605 21. Zaheer, R. *et al.* Impact of sequencing depth on the characterization of the microbiome
606 and resistome. *Sci. Rep.* **8**, 5890 (2018).

607 22. Ellington, M. J. *et al.* The role of whole genome sequencing in antimicrobial susceptibility
608 testing of bacteria: report from the EUCAST Subcommittee. *Clin. Microbiol. Infect.* **23**, 2–
609 22 (2017).

610 23. McDermott, P. F. *et al.* Whole-Genome Sequencing for Detecting Antimicrobial
611 Resistance in Nontyphoidal *Salmonella*. *Antimicrob. Agents Chemother.* **60**, 5515–5520
612 (2016).

613 24. Rodriguez-R, L. M. & Konstantinidis, K. T. Estimating coverage in metagenomic data sets
614 and why it matters. *ISME J.* **8**, 2349–2351 (2014).

615 25. Hillmann, B. *et al.* Evaluating the Information Content of Shallow Shotgun
616 Metagenomics. *mSystems* **3**, e00069-18 (2018).

617 26. Taur, Y. *et al.* Intestinal domination and the risk of bacteremia in patients undergoing
618 allogeneic hematopoietic stem cell transplantation. *Clin. Infect. Dis.* **55**, 905–914 (2012).

619 27. Stewart, C. J. *et al.* Longitudinal development of the gut microbiome and metabolome in
620 preterm neonates with late onset sepsis and healthy controls. *Microbiome* **5**, 75 (2017).

621 28. Zhai, B. *et al.* High-resolution mycobiota analysis reveals dynamic intestinal translocation
622 preceding invasive candidiasis. *Nat. Med.* **26**, 59–64 (2020).

623 29. Tamburini, F. B. *et al.* Precision identification of diverse bloodstream pathogens in the
624 gut microbiome. *Nat. Med.* **24**, 1809–1814 (2018).

625 30. Magruder, M. *et al.* Gut uropathogen abundance is a risk factor for development of
626 bacteriuria and urinary tract infection. *Nat. Commun.* **10**, 5521 (2019).

627 31. Clausen, P. T. L. C., Zankari, E., Aarestrup, F. M. & Lund, O. Benchmarking of methods for
628 identification of antimicrobial resistance genes in bacterial whole genome data. *J.*
629 *Antimicrob. Chemother.* **71**, 2484–2488 (2016).

630 32. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods*
631 **9**, 357–359 (2012).

632 33. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler
633 transform. *Bioinformatics* **25**, 1754–1760 (2009).

634 34. Inouye, M. *et al.* SRST2: Rapid genomic surveillance for public health and hospital

635 microbiology labs. *Genome Med.* **6**, 90 (2014).

636 35. Lanza, V. F. *et al.* In-depth resistome analysis by targeted metagenomics. *Microbiome* **6**,
637 11 (2018).

638 36. Moran, R. A., Anantham, S., Holt, K. E. & Hall, R. M. Prediction of antibiotic resistance
639 from antibiotic resistance genes detected in antibiotic-resistant commensal *Escherichia*
640 *coli* using PCR or WGS. *J. Antimicrob. Chemother.* **72**, 700–704 (2017).

641 37. Mahfouz, N., Ferreira, I., Beisken, S., von Haeseler, A. & Posch, A. E. Large-scale
642 assessment of antimicrobial resistance marker databases for genetic phenotype
643 prediction: A systematic review. *J. Antimicrob. Chemother.* **75**, 3099–3108 (2020).

644 38. Zankari, E. *et al.* Identification of acquired antimicrobial resistance genes. *J. Antimicrob.*
645 *Chemother.* **67**, 2640–2644 (2012).

646 39. Goldstein, S., Beka, L., Graf, J. & Klassen, J. L. Evaluation of strategies for the assembly of
647 diverse bacterial genomes using MinION long-read sequencing. *BMC Genomics* **20**, 23
648 (2019).

649 40. Su, M., Satola, S. W. & Read, T. D. Genome-Based Prediction of Bacterial Antibiotic
650 Resistance. *J. Clin. Microbiol.* **57**, e01405-18 (2019).

651 41. Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering, reconstruction and
652 typing of plasmids from draft assemblies. *Microb. Genomics* **4** (2018).

653 doi:10.6084/m9.figshare.6177188

654 42. Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive
655 antibiotic resistance database. *Nucleic Acids Res.* **48**, D517–D525 (2019).

656 43. Pereira-Marques, J. *et al.* Impact of Host DNA and Sequencing Depth on the Taxonomic

657 Resolution of Whole Metagenome Sequencing for Microbiome Analysis. *Front. Microbiol.*

658 **10**, 1277 (2019).

659 44. Cho, M. Y. *et al.* Two-target quantitative PCR to predict library composition for shallow

660 shotgun sequencing. *bioRxiv* 1–14 (2020). doi:10.1101/2020.09.21.304006

661 45. Petrof, E. O. *et al.* Stool substitute transplant therapy for the eradication of Clostridium

662 difficile infection: ‘RePOOPulating’ the gut. *Microbiome* **1**, 3 (2013).

663 46. Li, H. Seqtk: Toolkit for processing sequences in FASTA/Q formats. *Github* (2012).

664 Available at: <https://github.com/lh3/seqtk>.

665 47. Andrews, S. FastQC: a quality control tool for high throughput sequence data. (2010).

666 48. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina

667 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).

668 49. Bankevich, A. *et al.* SPAdes: A New Genome Assembly Algorithm and Its Applications to

669 Single-Cell Sequencing. *J. Comput. Biol.* **19**, 455–477 (2012).

670 50. Nurk, S., Meleshko, D., Korobeynikov, A. & Pevzner, P. A. MetaSPAdes: A new versatile

671 metagenomic assembler. *Genome Res.* **27**, 824–834 (2017).

672 51. Truong, D. T. *et al.* MetaPhlAn2 for enhanced metagenomic taxonomic profiling. *Nat.*

673 *Methods* **12**, 902–903 (2015).

674 52. Camacho, C. *et al.* BLAST+: Architecture and applications. *BMC Bioinformatics* **10**, 421

675 (2009).

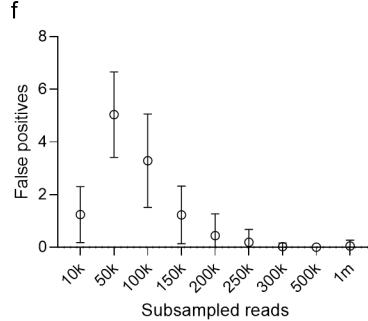
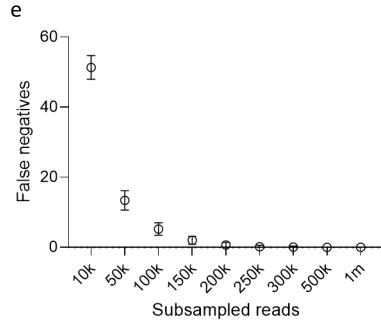
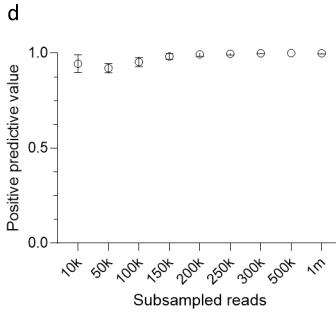
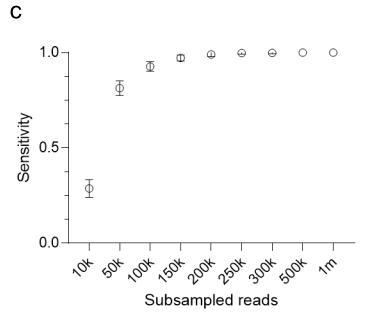
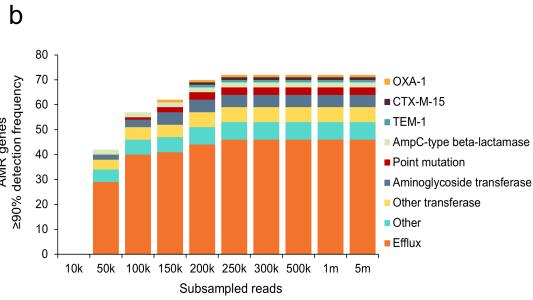
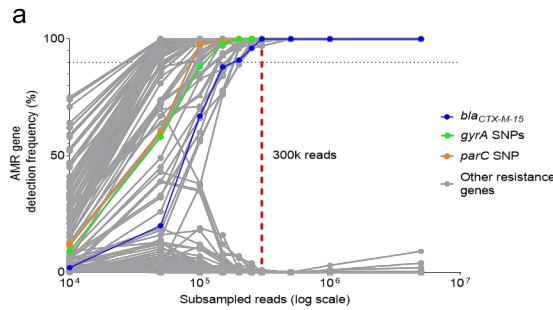
676 53. Hyatt, D. *et al.* Prodigal: Prokaryotic gene recognition and translation initiation site

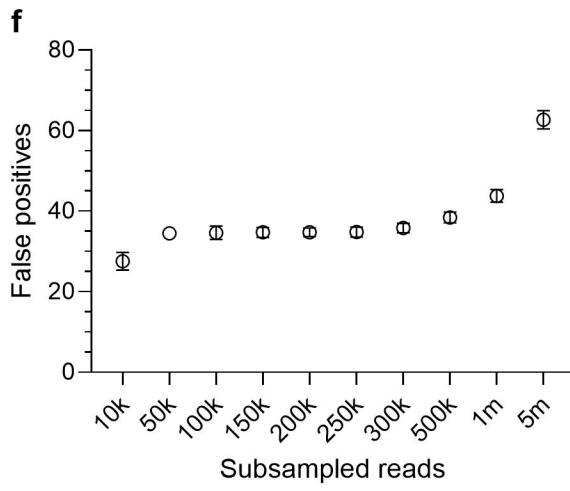
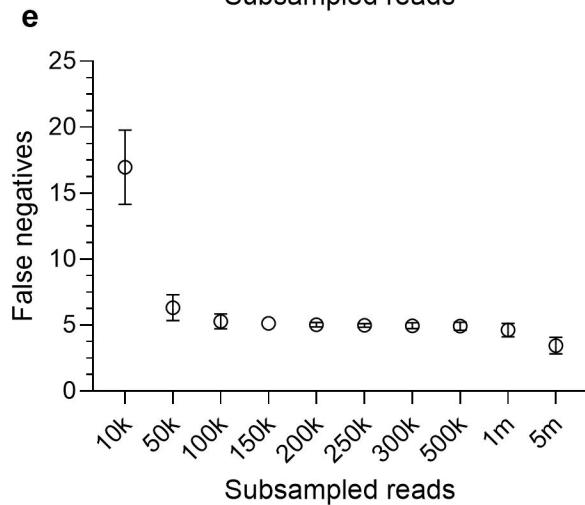
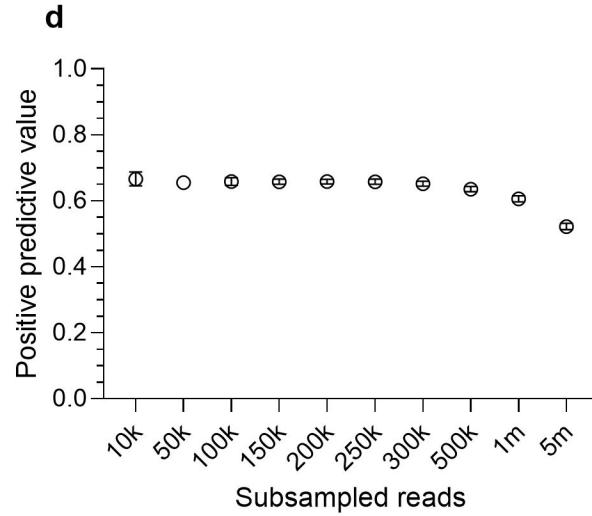
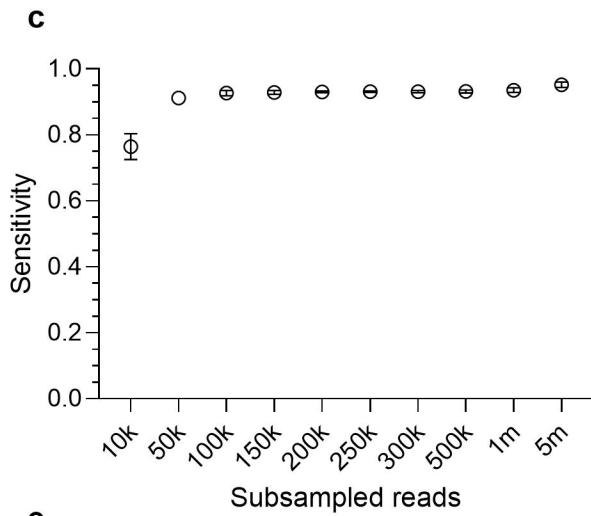
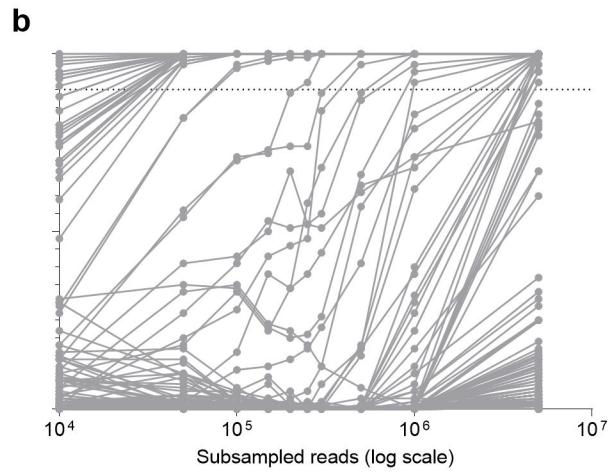
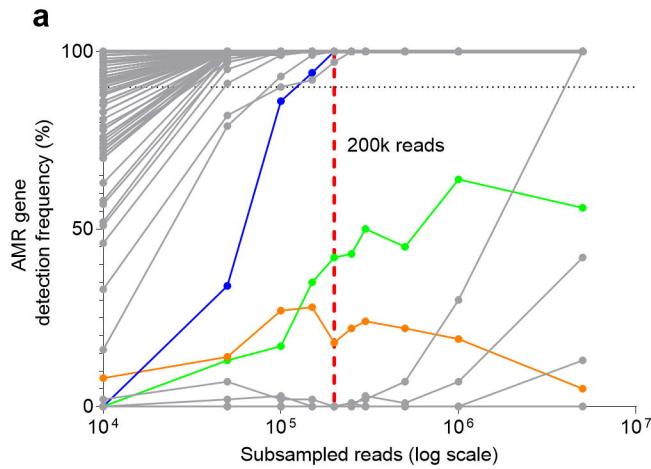
677 identification. *BMC Bioinformatics* **11**, 119 (2010).

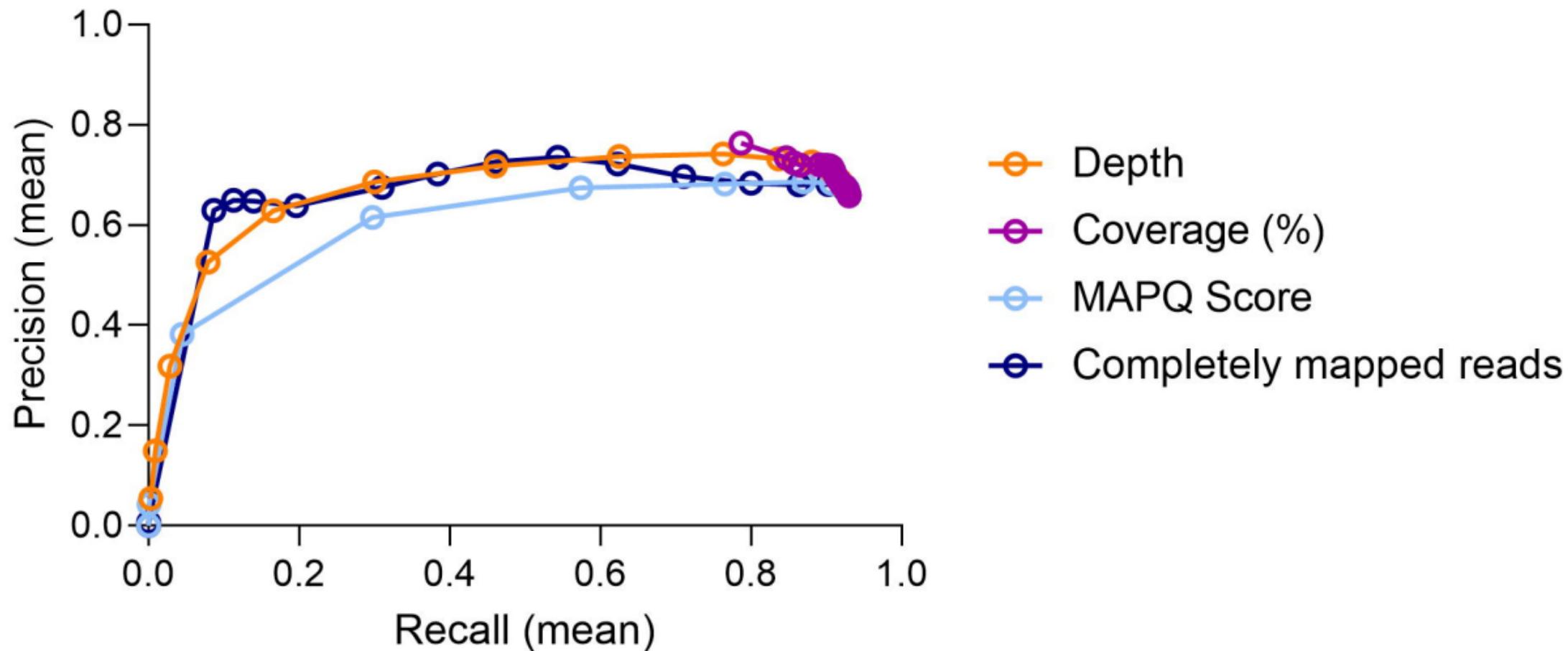
678 54. Lander, E. S. & Waterman, M. S. Genomic mapping by fingerprinting random clones: A

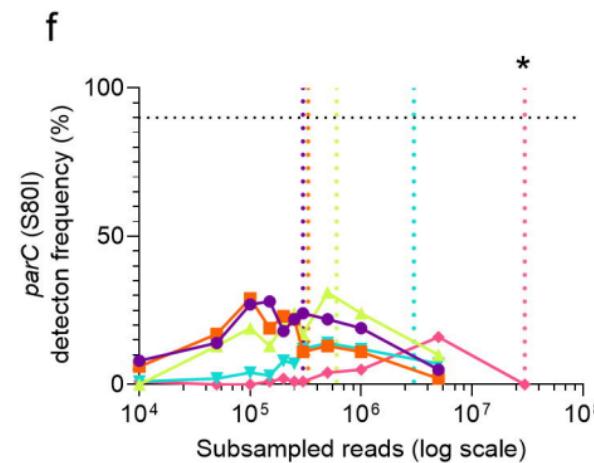
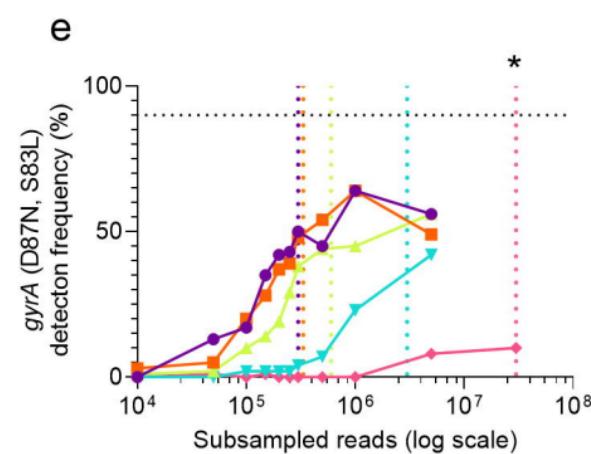
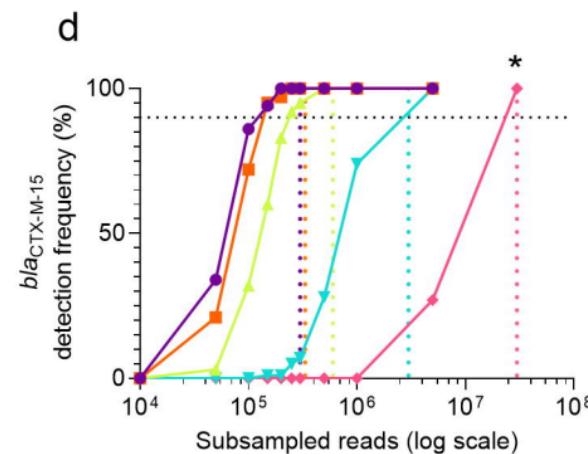
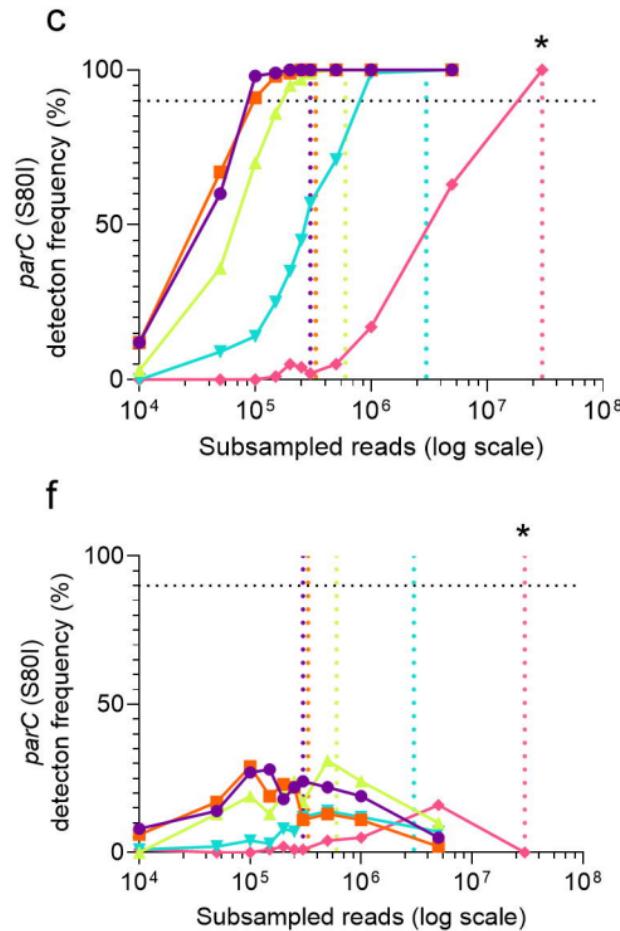
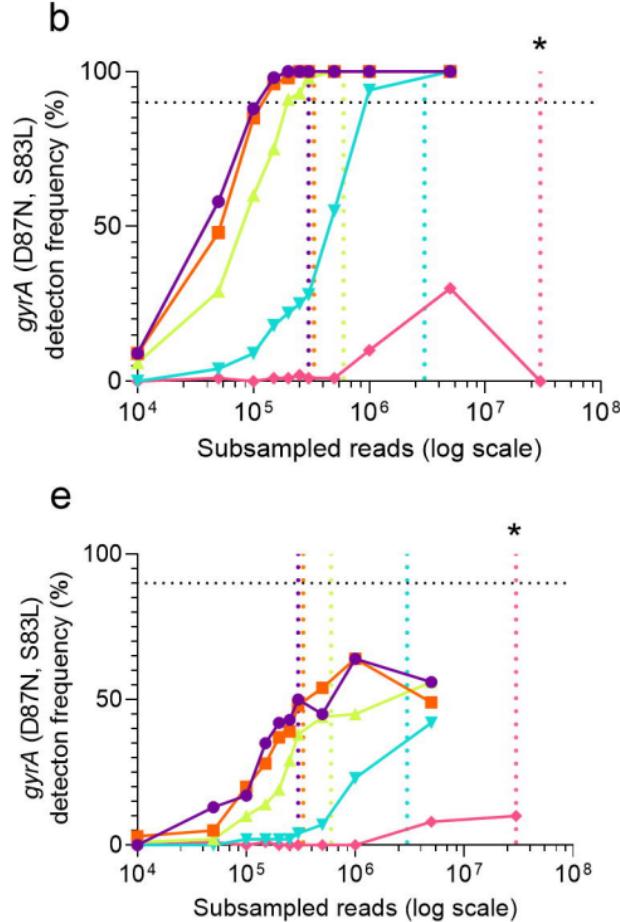
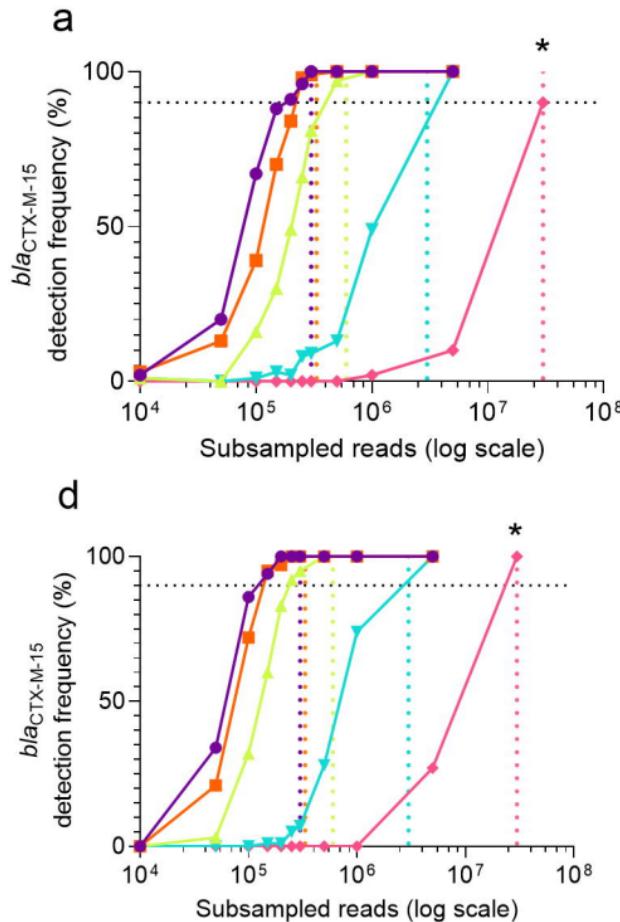
679 mathematical analysis. *Genomics* **2**, 231–239 (1988).

680









E. coli ST38 relative abundance

- 100% (Isolate)
- 90%
- 50%
- 10%
- 1%

