1    **Title: Speed, accuracy, sensitivity and quality control choices for detecting clinically**
2    **relevant microbes in whole blood from patients**
3
4    **Short title:** Detecting pathogens in clinically relevant samples
5

6    **Authors**: James Thornton Jr.[2]*, George S. Watts[1]*, Ken Youens-Clark[2], Lee D. Cranmer[3], and
7    Bonnie L. Hurwitz[2,4]†
8
9    **Affiliations:**
10   [1]The University of Arizona Cancer Center and Department of Pharmacology, The University of
11   Arizona, Tucson, AZ, USA
12   [2]Department of Biosystems Engineering, The University of Arizona, Tucson, AZ, USA
13   [3]Department of Medicine at the University of Washington, Fred Hutchinson Cancer Research
14   Center, and Seattle Cancer Care Alliance, Seattle, WA, USA
15   [4]BIO5 Institute, The University of Arizona, Tucson, AZ, USA
16   *These authors contributed equally to this work.
17   † To whom correspondence should be addressed
18
19   **Corresponding author**: bhurwitz@email.arizona.edu

1

20    **ABSTRACT**

21

22         Infections are a serious health concern worldwide, particularly in vulnerable populations

23    such as the immunocompromised, elderly, and young. Advances in metagenomic sequencing

24    availability, speed, and decreased cost offer the opportunity to supplement or replace culture-

25    based identification of pathogens with DNA sequence-based diagnostics. Adopting metagenomic

26    analysis for clinical use requires that all aspects of the pipeline are optimized and tested,

27    including data analysis. We tested the accuracy, sensitivity, and resource requirements of

28    Centrifuge within the context of clinically relevant bacteria. Binary mixtures of bacteria showed

29    Centrifuge reliably identified organisms down to 0.1% relative abundance. A staggered mock

30    bacterial community showed Centrifuge outperformed CLARK while requiring less computing

31    resources. Shotgun metagenomes obtained from whole blood in three febrile neutropenia patients

32    showed Centrifuge could identify both bacteria and viruses as part of a culture-free workflow.

33    Finally, Centrifuge results changed minimally by eliminating time-consuming read quality

34    control and host screening steps.

35
36
37    **AUTHOR SUMMARY**

38

39         Immunocompromised patients, such as those with febrile neutropenia (FN), are

40    susceptible to infections, yet cultures fail to identify causative organisms ~80% of the time.

41    High-throughput metagenomic sequencing offers a promising approach for identifying pathogens

42    in clinical samples. Mining through metagenomes can be difficult given the volume of reads,

43    overwhelming human contamination, and lack of well-defined bioinformatics methods. The goal

44    of our study was to assess Centrifuge, a leading tool for the identification and quantitation of

45    microbes, and provide a streamlined bioinformatics workflow real-word data from FN patient

1

46 blood samples. To ensure the accuracy of the workflow we carefully examined each step using

47 known bacterial mixtures that varied by genetic distance and abundance. We show that

48 Centrifuge reliably identifies microbes present at just 1% relative abundance and requires

49 substantially less computer time and resource than CLARK. Moreover, we found that Centrifuge

50 results changed minimally by quality control and host-screening allowing for further reduction in

51 compute time. Next, we leveraged Centrifuge to identify viruses and bacteria in blood draws for

52 three FN patients, and confirmed suspected pathogens using genome coverage plots. We

53 developed a web-based tool in iMicrobe and detailed protocols to promote re-use.

54

**INTRODUCTION**

The current gold standard for clinical diagnosis of infections relies on isolating organisms by culture-based methods followed by identification and drug resistance testing. Methods for identifying pathogens that rely on culture have several drawbacks including fastidious bacteria, the time required for growth in culture, and the difficulty targeting viruses, fungi, and parasites. Identifying pathogens directly from biological samples by DNA sequencing can overcome the above limitations of culture and may improve the rate and speed of diagnosis. For these reasons, metagenomic shotgun sequencing of pathogens has been referred to as the holy grail of infection diagnosis (Ecker et al., 2010). While culturing samples is the current standard for infection diagnosis, it can have a high failure rate in some scenarios. For example, a study examined the problem of culture-based diagnosis of infection in febrile neutropenia and found that only ~16% (609 of 3,756) febrile neutropenia patients were culture positive (van Walraven & Wong, 2014). Also, the hazard ratio of dying was nearly four-fold higher in culture-negative patients than for patients where no culture was taken (presumably due to lack of fever), indicating the high cost in lives when cultures fail. Therefore, we seek to apply metagenomic sequencing to overcome the low rate and time delay of culture-based diagnostic methods in clinical settings such as febrile neutropenia.

The potential of metagenomic shotgun sequencing has been demonstrated in a broad range of infection scenarios including: leptospirosis (Wilson et al., 2014), nosocomial transmission of a drug-resistant bacteria (Snitkin et al., 2012), foodborne illness (Ashton et al., 2015), and infectious disease outbreaks (Quick et al., 2016). Despite successes using metagenomic shotgun sequencing to identify pathogens, routine application in clinical settings will require accurate, efficient classification, with minimized sample contamination. For

3

79  example, while a small group of studies have reported on high-throughput metagenomic

80  sequencing for identifying pathogens from immunocompromised patients where samples were

81  not enriched for microbes, resulting in less than 1% of reads being pathogen-specific (Naccache

82  et al., 2014; Parize et al., 2017) and dramatically reducing the diagnostic possibilities from the

83  data (Frey et al., 2014). To begin addressing these inefficiencies, we developed an approach to

84  increase the proportion of pathogen-derived reads in samples and applied it to the patient

85  samples reported here.

86       On the data analysis side, there are no standards for analysis of metagenomic data

87  obtained from clinical samples; however, there have been recent innovations in taxonomic

88  classification algorithms that make it possible to quantify microbial species directly from reads

89  in metagenomic datasets rapidly. These algorithms use two main approaches to assign reads to

90  species in a reference database including: (1) a mapping approach using a Burrows-Wheeler

91  transform (Li & Durbin, 2009; M. Burrows, 1994) used by Centrifuge (Kim, Song, Breitwieser,

92  & Salzberg, 2016) or (2) a pseudo-alignment approach based on discriminating k-mers used by

93  CLARK (Ounit, Wanamaker, Close, & Lonardi, 2015a). These algorithms outperform local

94  alignment methods concerning both speed and capacity and can, therefore, better handle the

95  number of reads in metagenomes (Bazinet & Cummings, 2012; Ounit, Wanamaker, Close, &

96  Lonardi, 2015b; Rosen, Reichenberger, & Rosenfeld, 2011; Wood & Salzberg, 2014). However,

97  comparisons between these algorithmic approaches to determine the accuracy of taxonomic

98  assignment in clinically relevant metagenomes are lacking.

99       Here we report the accuracy and sensitivity of Centrifuge utilizing defined clinically

100  relevant samples, compare its performance to CLARK, and finally analyze datasets obtained

101  from patients following depletion of human cells to enrich for pathogen DNA. Lastly, we test the
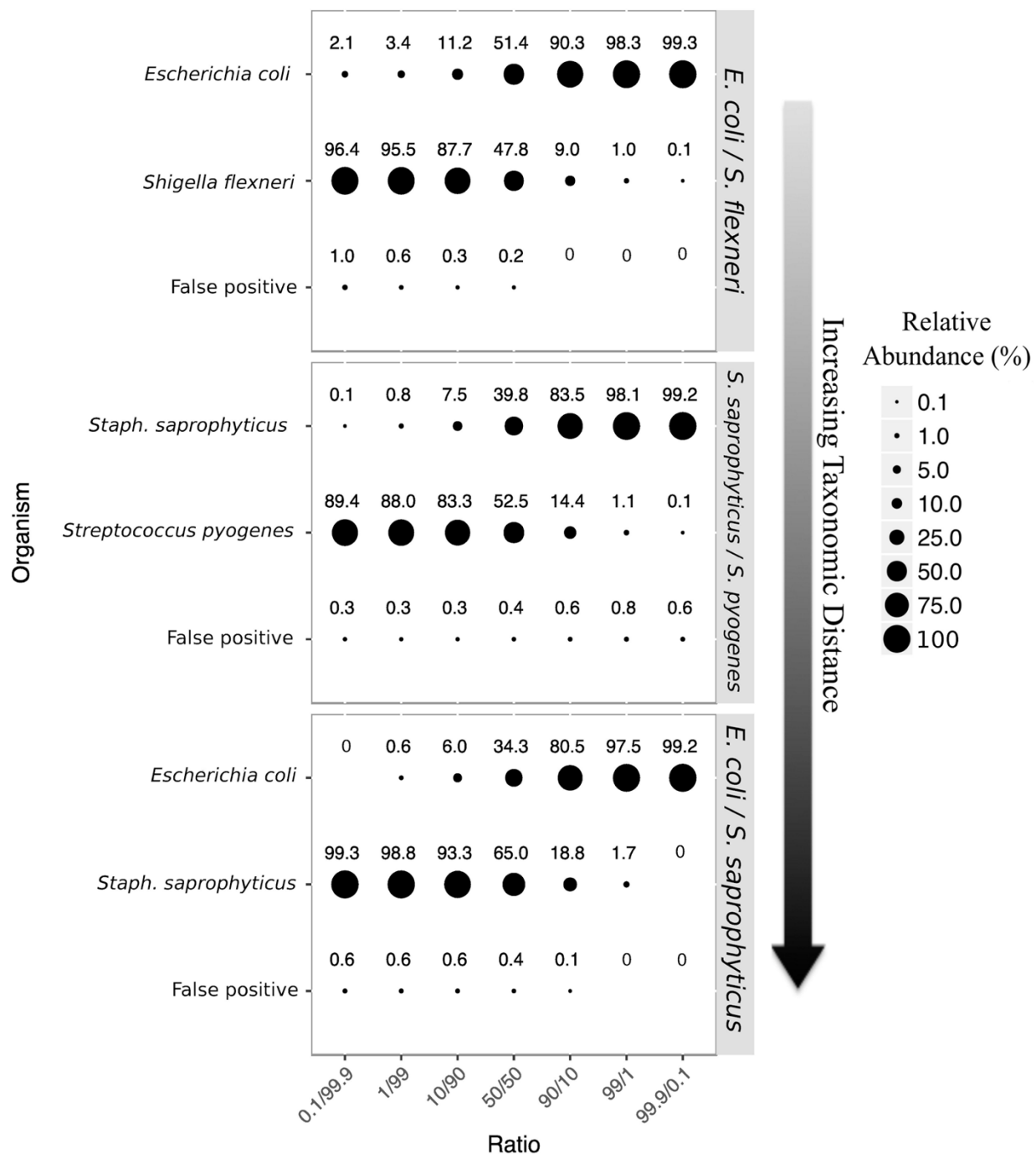
4

102   effect of excluding quality control and host-screening by alignment on the classification of reads

103   by Centrifuge. This work provides a foundation for analysis of metagenomic data from clinical

104   samples enriched for pathogens which use open-source software, requires a minimal

105   computational resource, and provides rapid and accurate identification of pathogens. Our

106   approach is freely available as web-based Apps in iMicrobe. Further, we provide the source code

107   in GitHub: https://github.com/hurwitzlab/Centrifuge_HPC under the GNU open source license.

108
109   **RESULTS**
110
111   **Centrifuge accuracy and sensitivity in controlled mixtures of bacteria**
112
113        Because closely related clinically important bacteria can have diametric clinical

114   consequences, (e.g., *E. coli* is a normal commensal while *S. flexneri* causes dysentery), we

115   sought to test Centrifuge's appropriateness as a tool for analyzing clinically relevant bacterial

116   sequence datasets. We tested the linearity and threshold for detection of Centrifuge using three

117   sets of bacterial mixtures, selected to represent taxonomic distances from phylum to genus-level.

118   We created dilution mixtures over a six-log range of relative abundance with each organism

119   ranging from 0.1% to 99.9% of the mixture (Figure 1). Centrifuge correctly identified all four

120   species in the mixtures and misidentified less than one percent of the reads in any of the 18

121   combinations sequenced (false positives, Figure 1). Centrifuge was sensitive to the lowest

122   relative abundance (0.1%) in four out of six opportunities, failing to detect the extremes in the *E.*

123   *coli/S. saprophyticus* mixture. Reads matching phage present in the mixtures were classified and

124   quantitated by Centrifuge separately from their host genomes. Because the phage relative

125   abundance estimates were not included with their host, the bacteria present were underestimated

126   so that the abundance estimates shown in Figure 1 do not add to 100%. The clearest example of

127   phage matches affecting taxon-assignment is in the mixture composed of 99.9% *S. pyogenes*

5

128     with an estimated relative *abundance of Streptococcus*-specific phage at 10.14%. Despite the

129     effect of phage matches, the coefficient of determination ($R^2$) for the three mixtures was 0.90 for

130     *E. coli*/*S. flexneri*, 0.99 for *S. saprophyticus*/*S. pyogenes*, and 0.96 for *E. coli*/*S. saprophyticus*.

131     Importantly, Centrifuge was able to discriminate between organisms as difficult to discriminate
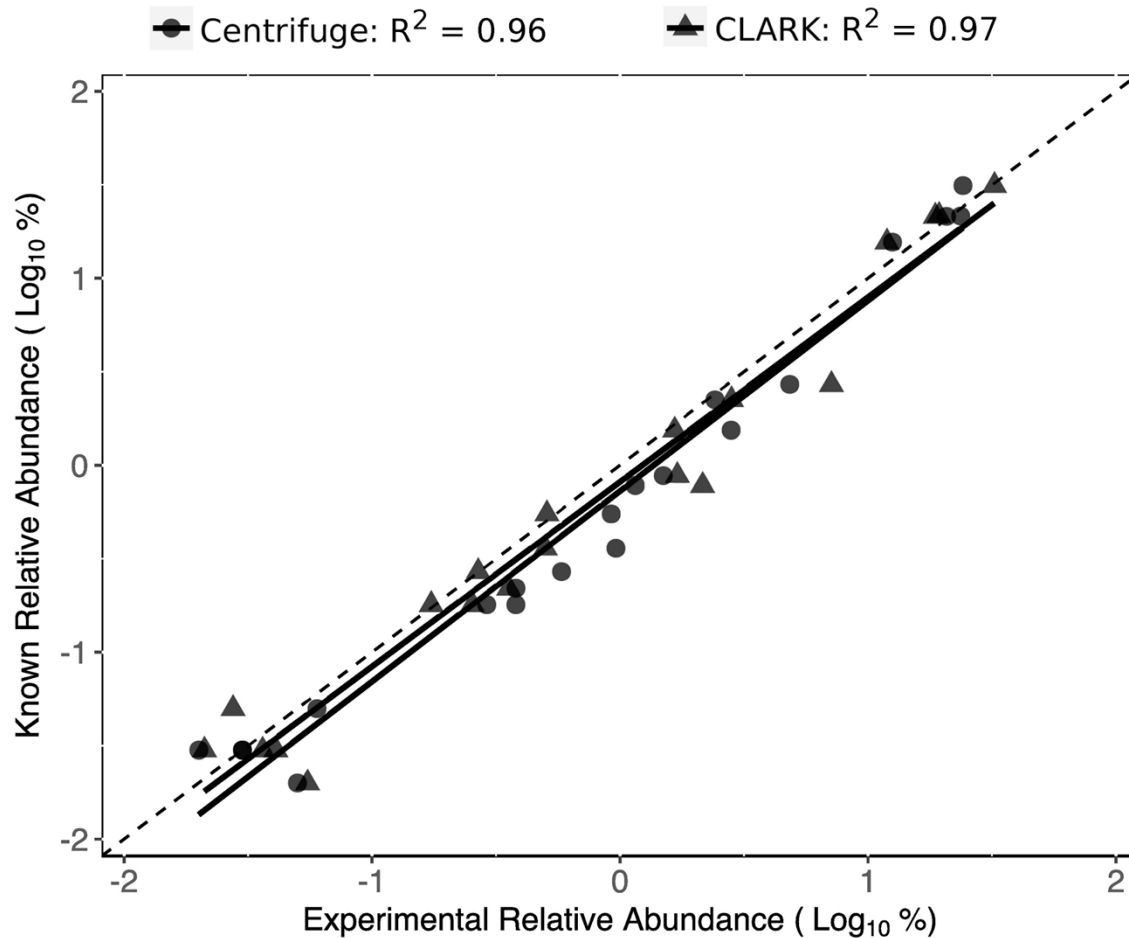
132     as *E. coli* and *S. flexneri*.

133

**Fig 1. Linearity and threshold for detection of binary mixtures of bacteria using Centrifuge**. The relative abundance of organisms calculated by Centrifuge is represented by circle size with actual values displayed above, values that are zero have no circle.

140  **Comparing the accuracy of Centrifuge and CLARK with a bacterial mock community**

141     Given Centrifuge's performance on the binary mixtures, it was next compared to a

142  leading algorithm of another class, CLARK with a more complex mock community of 20

143  bacteria present in varying relative abundances. Both CLARK and Centrifuge identified the 20

144  known bacterial species in the mock community; however, CLARK reported five false positives

145  (two *Shigella sp.*, two *Staphylococcus sp.* and *Corynebacterium pseudotuberculosis*) that were

146  not present in the mock community. In contrast to CLARK, Centrifuge did not produce any false

147  positives. To compare the two algorithms (Centrifuge and CLARK), we graphed the relative

148  abundance of 20 organisms in a mock community against their known abundance and calculated

149  R2 values (Figure 2). Centrifuge and CLARK had nearly identical $R^2$ values of 0.98 and 0.97

150  respectively. Overall, both tools tended to overestimate relative abundance values, especially the

151  lowest abundances: most estimated abundances fell below the perfect fit represented by the

152  dotted line in Figure 2. Importantly, both algorithms were able to identify the presence of all four

153  organisms in the mock community with relative abundances of 0.01%.

154

8

**Fig 2. Centrifuge and CLARK relative abundance estimates versus expected for a mock community of 20 bacteria.** Relative abundances estimated by CLARK and Centrifuge graphed against the expected values. The black dotted line represents perfect correlation with known relative abundances. The trendlines for CLARK and Centrifuge are shown in solid black lines.

**Centrifuge requires less computational resources than CLARK**

While CLARK had nearly identical accuracy in relative abundance estimates as Centrifuge (despite five positive identifications), there was a striking difference between the two classification algorithms in the computation resources and time required to analyze the data.

9

166 Relative to CLARK, Centrifuge required less than a tenth of the memory and a quarter of the

167 runtime, while using half the number of central processing units (Table 1).

168

169 **Table 1. Comparison of computational resources required by Centrifuge and CLARK to**

170 **analyze the bacterial mock community dataset.** CPU, central processing unit; GB, gigabyte;

171 RAM, random access memory.

| Program | number of CPUs | RAM (GB) | Runtime (hr:min:sec) |
|---|---|---|---|
| Centrifuge | 12 | 23 | 0:07:40 |
| CLARK | 28 | 297 | 0:38:40 |

172

173 **Identification of pathogens in whole blood from febrile neutropenia patients.**

174 Pathogens were enriched using a simple sample preparation method from whole blood

175 samples drawn from three patients with febrile neutropenia, and the resulting metagenomic DNA

176 sequenced. Table 2 shows the starting number of raw reads and the percent passing through each

177 step from quality control, to host-screening by alignment, and finally Centrifuge analysis. The

178 reads classified by Centrifuge identified three likely pathogens: Pseudomonas fluorescens with a

179 relative abundance of 50.7% in patient 1, Human parvovirus with a relative abundance of 99.8%

180 in patient 2, and Torque teno virus in patient 3 with a relative abundance of 62.8% (Figure 3).

181 Comparing the percentages shown in Table 2 with the relative abundances calculated by

182 Centrifuge for these organisms showed how the small genome sizes of the two viruses gave their

183 genomes more weight in the relative abundance estimates. For example, Torque Teno Virus had

184 an abundance estimate of 72.8% though only 9.4% of the total post-quality control reads mapped

185 to this organism.

186          Blood culture results for all three patients were negative, at the time of sample collection

187    and in two subsequent blood cultures of each patient. Thus, the sequencing results were not

188    compared to culture, the current gold standard. However, patient two did have a positive PCR

189    test for human parvovirus in the month before and after the research sample was obtained,

190    corroborating the results obtained with Centrifuge. Additional corroboration of the results comes

191    from analysis of 12 samples obtained from two healthy volunteers over a six-week period in

192    which none of the likely pathogens seen in the febrile neutropenia patients was observed (data

193    not shown). While *Pseudomonas fluorescens* has been reported as a false positive in other

194    studies, the fact that it did not appear in the healthy volunteer samples and is known to infect

195    immunocompromised individuals (Wong et. al., 2011) suggests that it is not an artifact in patient

196    1 (. We also identified human endogenous retrovirus K113 and *Cutibacterium acnes* (also known

197    as *Propionibacterium acnes*) in patients 1 and 3, however these organisms were deemed to be

198    contaminants: the virus is endogenous, *C. acnes* is a common contaminant of blood samples

199    (Mollerup et al., 2016; Parize et al., 2017; Park et al., 2011), and both were present in the normal

200    samples collected over 6 weeks.

201

202

203

204

205

206

207

208

209

210

211

212

213

214

**Table 2.** Read counts following each step of the Centrifuge analysis of febrile neutropenia

datasets. QC, quality control.

| Pt | Raw Reads[a] | Post QC[b] | Human (%)[c] | Centrifuge | | |
|---|---|---|---|---|---|---|
| | | | | Unmapped (%)[d] | Classified (%)[e] | Unknown (%)[f] |
| 1 | 3,497,123 | 61.9 | 57.3 | 42.7 | 70.2 | 29.8 |
| 2 | 13,000,518 | 43.9 | 41.3 | 58.7 | 34.8 | 64.2 |
| 3 | 18,839,275 | 43.4 | 79.1 | 20.9 | 45.4 | 54.6 |

217   [a] Total reads generated from the sample.

218   [b] Percent of reads remaining after quality control.

219   [c] Percent of Post-QC reads that mapped to the human genome.

220   [d] Percent of Post-QC reads that did not map to the human genome.

221   [e] Percent of unmapped reads that were assigned a taxonomic classification by Centrifuge.

222   [f] Percentage of unmapped reads that were not assigned a taxonomic classification by Centrifuge.

223
224

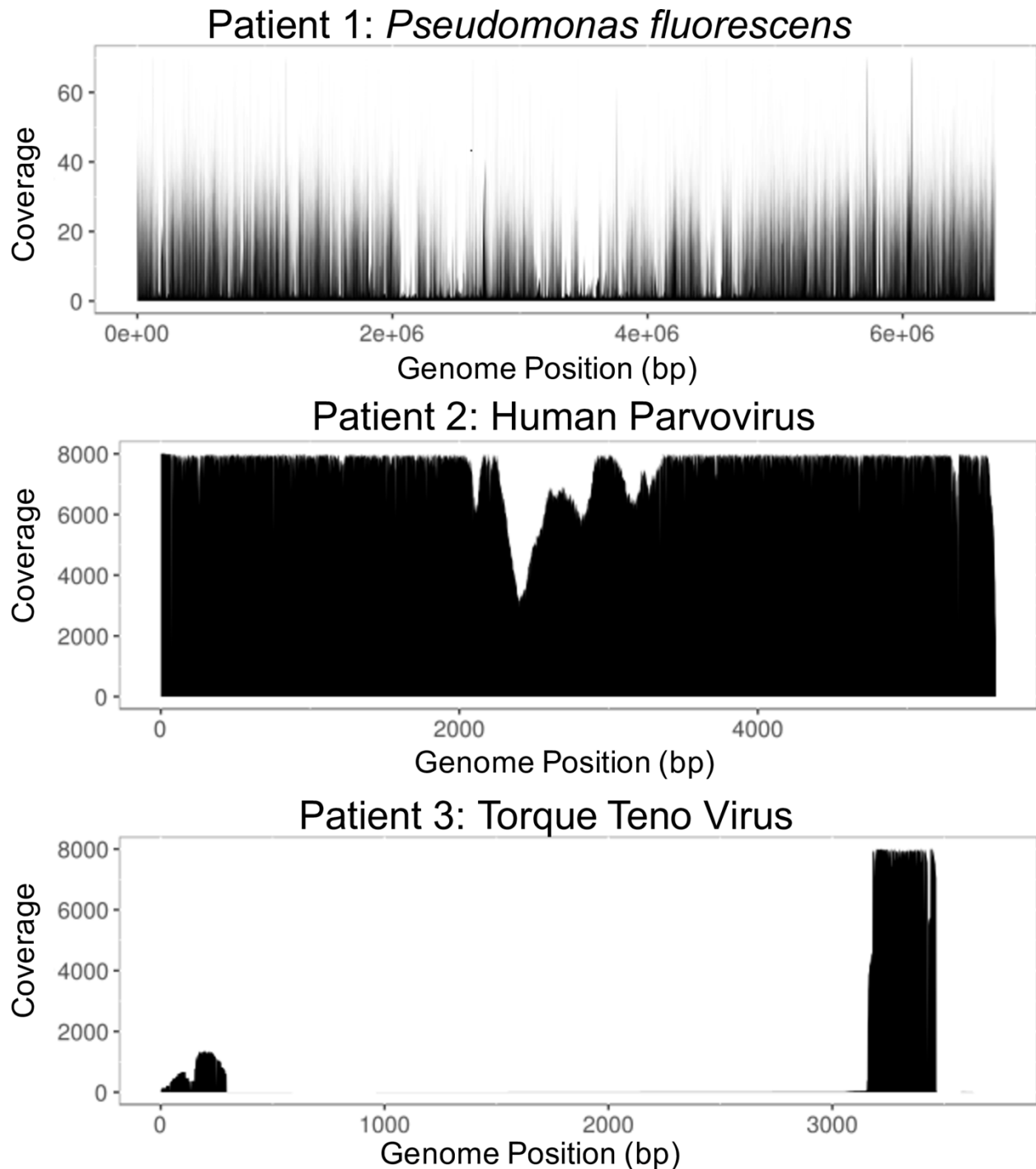**Fig 3. Identification and relative abundance of pathogens in febrile neutropenia samples**.

Circle size indicates the relative abundance of the respective organism, and actual abundance

values are next to the circles. Organisms deemed endogenous or common contaminants are

separated from the presumed pathogens by the horizontal line.


**Genome coverage of suspected pathogens in febrile neutropenic patients**

Reads from the three febrile neutropenia samples were aligned to the respective reference

genomes of the suspected pathogens to determine average depth of coverage (Figure 4). When

patient 1 reads were aligned to the *Pseudomonas fluorescens* genome, the average coverage was

7.0. Patient 2 reads aligned to the Human Parvovirus B19 genome showed average coverage of

5,180. Finally, patient 3 reads aligned to the Torque Teno Virus (TTV) genome showed high

coverage (~8,000) for a ~500 base pair region of the genome.

13

240



241
242
243 **Fig 4. Genome coverage of suspected pathogens identified in febrile neutropenia patients**

244 **by Centrifuge**. Coverage by each base is graphed relative to the position in the three respective
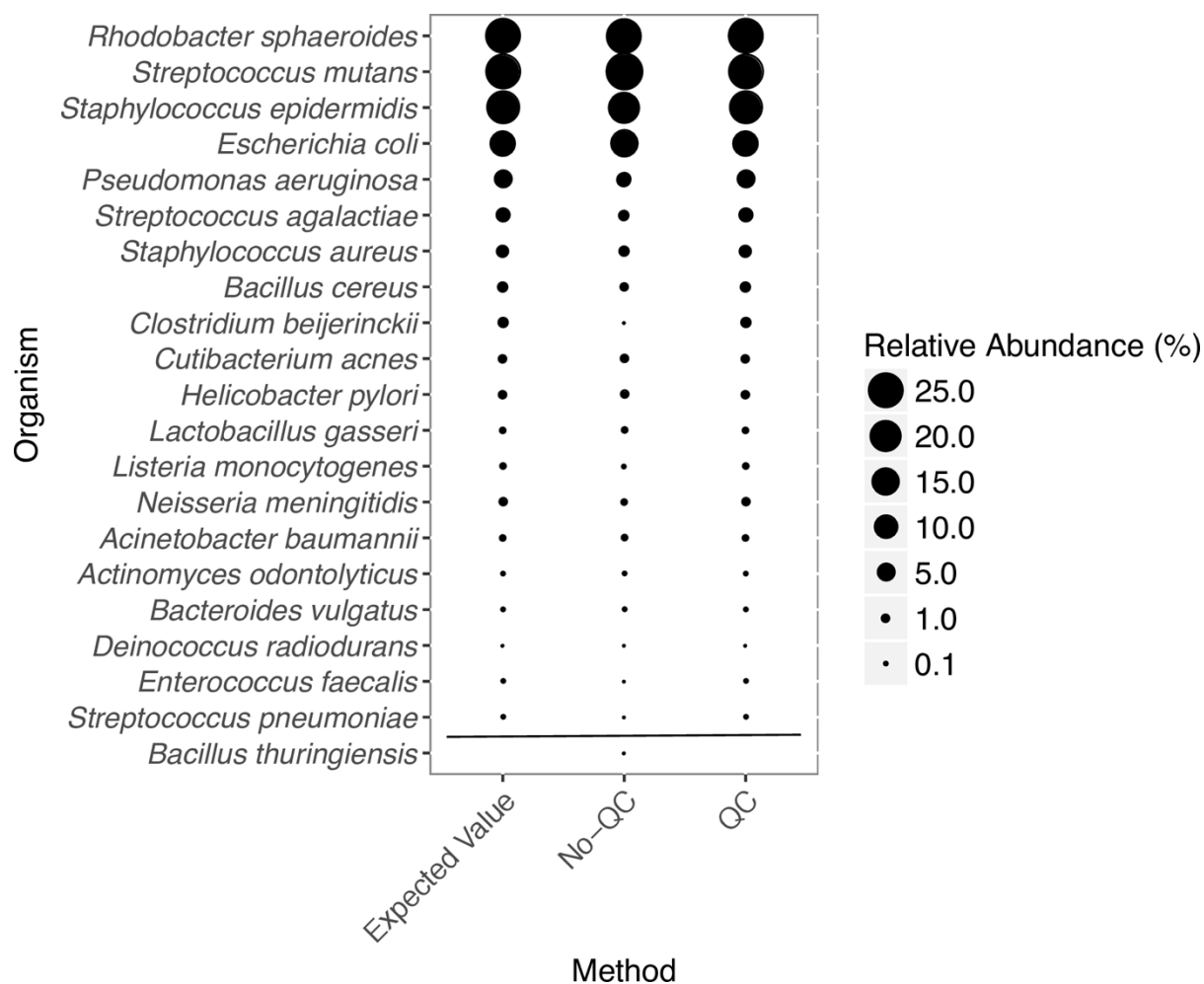
245 genomes of likely pathogens identified in three febrile neutropenia patients.

246

**Effect of quality controlling reads on computation time and Centrifuge's accuracy**

Sequencing reads are typically subjected to a series of quality control steps including

trimming low-quality bases from reads, removing short reads, deduplication, and trimming ends

with unbalanced nucleotide composition before downstream applications (e.g., variant calling, or

sequence assembly). When quality control steps were performed before the Centrifuge analyses

in Figures 2 and 3, they accounted for approximately half the compute time required to achieve

results (data not shown). The fact that quality controls steps accounted for so much of the

compute time, led to the question of what effect quality control had on the taxonomic

classifications and relative abundance estimates made by Centrifuge. To answer this question,

the mock bacterial community data was analyzed in Centrifuge with and without quality

controlling the reads first. Results showed only one difference in taxonomic classification: a false

positive (*Bacillus thuringiensis*) was identified with a relative abundance of 2.9% without quality

control (Figure 5). Linear regression of the measured versus expected relative abundances

showed that the $R^2$ with quality control was 0.97 and without quality control was 0.97, further

demonstrating how little effect there was on the Centrifuge results.

262

**Fig 5. Bacterial mock community taxonomic identification and relative abundance by Centrifuge with and without quality control of the input sequence reads.** Organisms are ranked by their relative abundance which is indicated by the size of the circle. The false positive (*Bacillus thuringiensis*) identified from reads without quality control (QC) is shown at the bottom.

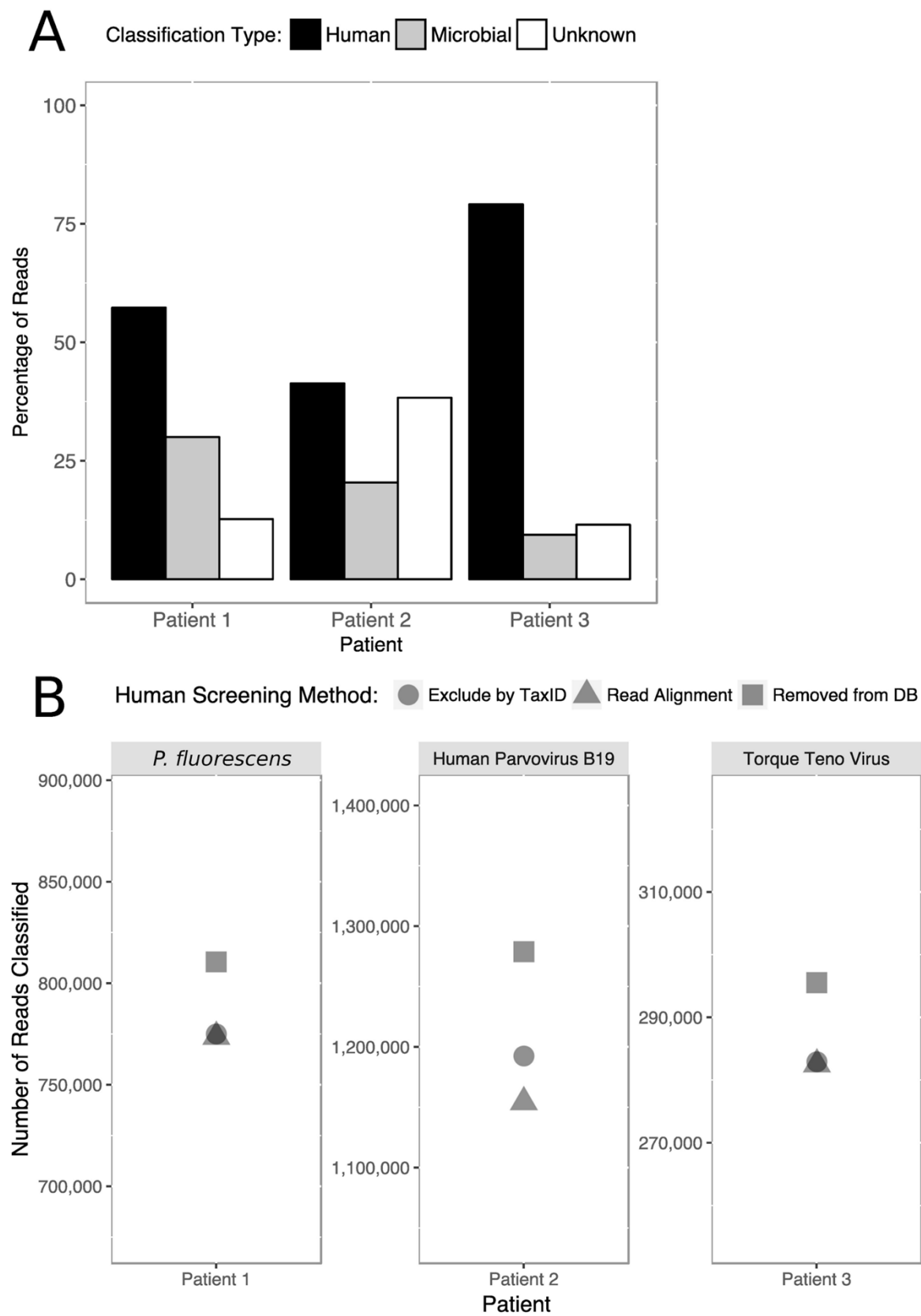**Host read removal by alignment versus in Centrifuge**

Host DNA contamination can contribute to a significant proportion, or even the vast majority, of reads in metagenomic datasets, and is often removed by mapping reads to the host genome (Schmieder & Edwards, 2011). In performing taxonomic classification of reads,

16

275  Centrifuge determines whether reads are of human origin (or other hosts), thus calling into

276  question the necessity of aligning reads to the host genome and removing them, before analysis.

277  Figure 6A shows the relative amount of reads that were classified as human, microbial, or

278  unknown when the datasets were analyzed by Centrifuge without removing reads by alignment

279  to the human genome before analysis. The relative proportion of host (human) reads in the data

280  agreed well with the proportions found by alignment (see Table 2). While the proportion of host

281  DNA was less than in prior studies, suggesting that the enrichment for pathogen DNA used in

282  this study was successful, a significant proportion of the reads were still human.

283       Having established that a significant proportion of the reads in the datasets were of host

284  origin by both alignment and Centrifuge, we compared three approaches for removing host reads

285  in the febrile neutropenia patient data. These methods include (1) alignment to the human

286  genome and removal of aligned reads from the dataset, (2) removing the human sequence from

287  the reference database, and (3) using the "exclude TaxID" function in Centrifuge to exclude

288  reads from classification whose best match was to the human genome. Overall, exclusion of the

289  human genome from the reference database resulted in the highest number of reads classified to

290  the presumed pathogens; however, the differences between the methods were relatively minor

291  (Figure 6B). Patient 3, with a presumed pathogen of Torque Teno virus, showed the least effect

292  on the number of reads classified, with less than a 411 read difference (<1%) between the

293  number of reads classified between the three methods. In contrast, patient 2, with a presumed

294  pathogen of Human Parvovirus B19, had 124,544 fewer reads classified (9.7%) when reads

295  removed by alignment relative to the removal of reads that match the human genome from the

296  database. Finally, patient 1, with a presumed pathogen of *P. fluorescens*, showed 26,836 fewer

17

297     reads classified (4.5%) when reads were removed by alignment relative to being present in the

298     database.

299
300

301 **Fig 6. Effect of mapping reads from patient samples against the human genome.** A) Percent

302 of reads classified as human, microbial, or unknown for each febrile neutropenia patient by

303 Centrifuge. B) Number of reads classified in each presumed pathogen following three strategies

304 for host screening: removal of the human genome from the reference database (Removed from

305 DB, squares), excluding the human TaxID in Centrifuge (Exclude by TaxID, circles), and

306 aligning against the human genome before analysis (Read Alignment, triangles).

307

308 **DISCUSSION**

309
310 **Centrifuge accuracy of identification and quantitation with known samples**

311 Immunocompromised patients, such as those with febrile neutropenia, are susceptible to

312 infections. The current standard for identifying pathogens from clinical samples when infection

313 is suspected can fail as much as ~80% of the time. Without diagnostic information, clinicians'

314 first response is empirical antibiotic therapy in the hope that the organism is bacterial and

315 covered by the antibiotic(s) given. Metagenomic sequencing of clinical samples offers an

316 approach that bypasses the issues of culture, however, mining the resulting metagenomic

317 sequence can be slow and error-prone given the volume of reads, host read contamination, and

318 lack of well-defined bioinformatics methods. The goal of our study was to assess Centrifuge, a

319 leading tool for identification and quantitation of metagenomic data, using clinically relevant

320 datasets to establish its accuracy in microbial/viral identification and abundance estimates with

321 an eye toward reducing compute time.

322 The first dataset used to assess Centrifuge was a series of binary bacterial mixtures

323 chosen for their phylogenetic distance and mixed so that each pair was combined across six logs

324 of relative abundance. Centrifuge was able to discriminate the most closely related pair of

20

325     bacteria, *E. coli* and *S. flexneri*, even when one of the organisms was present as 0.1% of the

326     mixture. As the proportion of *E. coli* decreased, the relative abundance estimate diverged from

327     expected, so that the *E. coli* estimate was 2.1% when *E. coli* was only 0.1% of the mixture. The

328     same inaccuracy did not occur as the *S. flexneri* relative abundance decreased to 0.1%,

329     suggesting Centrifuge misidentified a portion of the *S. flexneri* genome as *E. coli* but not the

330     other way around. The difficulty classifying *S. flexneri* suggested by the fact that the false

331     positive rate increased from 0% to 1%, the highest measured, as *S. flexneri* relative abundance

332     increased. One likely cause for more relative matches to *E. coli* than *S. flexneri* is that *E. coli*

333     strains and isolates represent the most substantial fraction of the Centrifuge reference database.

334     False positive identification of *E. coli* using metagenomic methods has been previously

335     observed. McIntyre et al. (2017) saw similar false positive identification of *E. coli* when using

336     metagenomic classifiers on negative control sequences not belonging to any known

337     organism(McIntyre et al., 2017). The researchers also speculated that the reason for the false

338     positives is the overrepresentation of *E. coli* sequences in their reference dataset. Although

339     Centrifuge uses a modified FM-index to condense closely related genomes, the total file size of

340     basepairs maintained (unique + shared based on $\geq$ 99% identity) exceeds the relative file size of

341     all other species (Kim et al., 2016) giving it a higher probability for matches. This result suggests

342     that Centrifuge dampens the effect of multiple strains and isolate genomes using the modified

343     FM-index, but the effect is still present for highly abundant strains.

344         Centrifuge appears to be capable of detecting organisms even when they are present in

345     minor abundance, regardless of the phylogenetic distances between them. Overall, Centrifuge

346     read abundances closely match the expected relative abundance of bacterial mixtures for closely

347     and distantly related species. Interestingly, phylogenetic distance did not predict the accuracy of

348     relative abundance estimates. A reasonable assumption would be that as phylogenetic distance

349     increases, the number of discriminatory k-mers increase to allow for better read classification by

350     Centrifuge. Instead, we observed high classification accuracy for the most closely related pair (*E.*

351     *coli*/*S. flexneri)* from the same family. Less accuracy for the next pair (*S. pyogenes*/*S.*

352     *saprophyticus*) where both organisms were gram-positive and from the same phylogenetic class.

353     The highest accuracy for the most distant pair (*E. coli*/*S. saprophyticus*) where one organism was

354     gram-negative and the other gram-positive and only shared phylogenetic kingdom. Interestingly,

355     *S. pyogenes* is closely related to many *Streptococcus* genomes which may have limited the

356     number of distinguishing k-mers to classify reads at the species rather than genus level (data not

357     shown).

358            We compared Centrifuge's performance against another leading k-mer based taxonomic

359     classifier, CLARK, in analyzing sequence data from a more complex community of 20 bacteria.

360     The mock community was also mixed in varying relative abundances as with the binary

361     mixtures, albeit, in a different range (~0.01-35%). Abundance calculations between the two

362     algorithms were nearly identical across the relative abundance range; however, the processing

363     time and computational resources for CLARK were greater (Table 1). Also, CLARK had a

364     propensity for false positives, whereas Centrifuge did not. On the other hand, Centrifuge's results

365     had to be processed to account for the strain and phage-specific data generated. Such processing

366     would be a necessary part of adoption in a clinical setting, but Centrifuge's lack of false positives

367     and speed suggests it may be a good starting point for such a tool.

368            **Centrifuge identification and relative abundance estimates**

369     Centrifuge is unique from other taxonomic classifiers in that it provides Expectation –

370     Maximization (EM) calculation to determine relative abundance, rather than just read

22

371    proportional classification. The EM calculation proves useful in determining relative abundance

372    between organisms in samples with varying genome sizes. We demonstrated the benefit of

373    calculating abundance using Centrifuge's EM algorithm in the analysis of the febrile neutropenia

374    blood samples from patients 2 and 3 where viral matches were significantly underrepresented

375    when using read proportional classifications.

376         One drawback for clinical pathogen identification is that Centrifuge separates strain-level

377    counts, splitting reads among closely related strains which required manually summing strain

378    level abundances for reporting. Future iterations of Centrifuge could address this issue re-

379    analyzing the data with a reduced reference set of genomes based on the first round of analysis or

380    a reduced reference database. Lastly, current reference databases do not account for all of the

381    extant microbial/viral diversity that may be present in patients. However, this issue is being

382    addressed over time with the exponential growth in the number of microbial draft genomes

383    available (Land et al., 2015).

384

385    **Genome coverage of presumptive pathogens identified in patient samples**

386    We examined genome coverage statistics with the assumption that the genomes of the pathogens

387    identified as the presumed cause of fever in the patients would be represented by consistent

388    coverage, whereas uneven coverage could indicate insufficient evidence of organism presence.

389    Parize et al. took a similar approach in which even distribution of contigs was used as part of the

390    criteria to decide if a sample was deemed positive (Parize et al., 2017). Interestingly, the Torque

391    Teno virus sequence found in patient 3 was observed to have high coverage of only a ~500 base

392    pair untranslated region of the genome. This highly conserved region has been suggested to be

393    critical for viral replication that may indicate an early replication event or the presence of

23

394    subviral particles, a characteristic that has previously observed in Torque Teno virus (de Villiers,

395    Borkosky, Kimmel, Gunst, & Fei, 2011). The evidence for sub-viral particles provided by the

396    coverage analysis is the first from an *in vivo* sample. Lastly, Torque Teno virus was identified in

397    a cancer patient undergoing bone marrow ablation in preparation for a hematopoietic stem cell

398    transplant as part of their cancer treatment. This finding highlights the possible value of the

399    metagenomic sequencing approach as Torque Teno virus has been investigated as a predictive

400    marker for post-transplant complications (Wohlfarth et al., 2018).

401

402    **Quality control of reads before Centrifuge analysis**

403    Although quality control of raw reads is imperative for variant calling and genome assembly and

404    can speed up downstream taxonomic and functional analyses by reducing the total number of

405    reads analyzed, it takes considerable computing time and resources. In this study, we observed

406    limited benefits of quality control regarding accurately identifying and quantifying the

407    abundance of the bacteria in the mock community. However, we did see an elimination of a

408    single false positive organism estimated at 2.3% relative abundance with quality control. Quality

409    controlling reads from the febrile neutropenia data revealed a bias toward removing viral reads

410    (Supplemental Table 1). Users of Centrifuge may want to weigh the limited benefits of quality

411    controlling their data before analysis in Centrifuge versus the bias toward the removal of viral

412    reads and time required.

413

414    **Host screening with Centrifuge**

415    Despite the substantial enrichment for microbial/viral DNA that we achieved in this study (20-

416    58% non-human reads, Table 2) as compared to prior studies (1% of reads) (Naccache et al.,

24

417    2014; Parize et al., 2017), a large proportion of reads were still identified as human. Screening

418    host reads by alignment to the genome before analysis by Centrifuge appears to be unnecessary

419    given Centrifuge's ability to classify reads to the host organism during analysis. For example, in

420    patient 2 we were able to identify Human Parvovirus B19 when we used the "exclude TaxID"

421    function for host screening. Because parvovirus virus integrated into the ancestral human

422    genome during evolution (Liu et al., 2011), many Human Parvovirus B19 reads identified

423    aligned to the human genome and were removed before analysis by Centrifuge. This method

424    caused the largest reduction in the number of reads classified as Human Parvovirus B19 relative

425    to the exclude TaxID method (Figure 6B).

426    In contrast, when the human genome was removed from the Centrifuge database, reads from the

427    human genome derived from the ancestrally integrated parvovirus would have been misclassified

428    as Human Parvovirus B19, with the effect that it could inflate the relative abundance estimate.

429    The "exclude TaxID" method appears to offer a balance between the other two methods: it

430    allows both endogenous host reads and actual organism reads to be appropriately classified while

431    saving the time and computational cost of aligning reads to a host organism before analysis.

432    Given that reference genomes can contain sequences of mixed origin due to horizontal gene

433    transfer, endogenous and integrated microbes/viruses, and prophage in bacterial genomes,

434    classifying reads to all available reference data and then utilizing exclude TaxID appears to be

435    the best compromise of speed and specificity for eliminating host reads from results.

436

437

438    **Conclusion**

439    In summary, our analyses suggest that Centrifuge, open-source software for fast taxonomic

440    classification, provides accurate quantification of clinically relevant organisms/viruses in

441    metagenomes using minimal compute time and resources. Centrifuge's ability to quickly assign

442    taxonomy to reads, accurately represent the abundance of organisms such as viruses, and

443    sidestep read quality control and host-screening make it a good candidate for classifying reads of

444    clinically relevant organisms. To this end, we have made Centrifuge and the bubble plot software

445    used in the study available as Apps in iMicrobe (http://imicrobe.us) for streamlined taxonomic

446    analysis by the public.

447
448    **Materials and Methods**
449
450    These methods have been deposited into protocols.io under DOI:

451    dx.doi.org/10.17504/protocols.io.wjdfci6

452
453    **Ethics Statement**

454    The Institutional Review Board at the University of Arizona (project #1505826794) approved the

455    human subjects research. Informed consent was obtained from febrile neutropenia patients.

456    Whole blood was collected from patients that developed febrile neutropenia during their

457    treatment at the University of Arizona Cancer Center. Data obtained from the first three patients

458    collected as part of a more extensive study were used here. All three patients were being treated

459    for leukemia or lymphoma at the time of their febrile neutropenia diagnosis.

460
461    **Binary mixtures of bacteria**

462        The binary mixtures were described previously (Watts et al., 2017). Briefly, four species

463    of bacteria were used to create three binary mixtures representing: (1) difficult to discriminate

464    species with divergent clinical impact (*Escherichia coli* versus *Shigella flexneri*); (2) Gram-

26

465    positive species (*Staphylococcus saprophyticus* versus *Streptococcus pyogenes*); and (3) Gram-

466    positive versus Gram-negative species (*E. coli* versus *S. saprophyticus*). DNA from the bacteria

467    were purchased from the American Type Culture Collection (Manassas, Va, USA) and mixed in

468    pairs so that each species represented 99.9, 99, 90, 50, 10, 1, and 0.1% of the total sample.

469    Samples were sequenced as described below, and the sequence data deposited to the NCBI

470    Sequence Read Archive under accessions: SRX3154186-SRX3154219 in project accession

471    PRJNA401033.

472

473    **Staggered mock bacterial community**

474         The mock bacterial community (BEI Resources, Manassas, VA, USA, National Institute

475    Allergy and Infectious Diseases, National Institutes of Health, as part of the Human Microbiome

476    Project: Genomic DNA from Microbial Mock Community B (Staggered, High Concentration),

477    v5.2H, for metagenomic shotgun sequencing, HM-277D) consisted of 20 bacterial species

478    created as part of the Human Microbiome Project with specific staggered 16S rRNA gene

479    abundances for each species. Using the 16S rRNA gene copy values, along with the known 16S

480    rRNA gene copy number in each species' genome, we calculated the number of genomes present

481    for each species to provide an expected value for comparison to the relative abundances

482    calculated by Centrifuge and CLARK from sequencing data. The mock community was

483    sequenced as described below, and sequence data deposited to the NCBI Sequence Read Archive

484    under accession: SRP115095 in project accession PRJNA397434.

485

486

487    **Febrile neutropenia patient blood samples**

488        Approximately five milliliters of whole blood were collected ($K_2$EDTA BD Vacutainer

489    tubes, catalog #367863 BD Biosciences, San Jose, CA, USA) when blood cultures were ordered

490    for each patient and transferred for processing within 2 hours of collection. Blood samples were

491    diluted with an equal volume of sterile phosphate buffered saline, layered on Ficoll-Paque (GE

492    HealthCare Life Sciences, Pittsburgh, PA, USA) and centrifuged for 20 minutes at 400 x g.

493    Plasma was carefully drawn off, sacrificing some yield to prevent drawing up monocytes, and

494    centrifuged three more times at 50, 100, and 150 x g for 5 minutes to further remove human

495    cells. The plasma was passed through a five-micron filter and finally centrifuged at 4000 x g.

496    DNA was isolated from any material sedimented during the final centrifugation with a UCP Pure

497    Pathogen kit (Qiagen Inc., Germantown, MD, USA). Isolated DNA was quantitated on a

498    NanoDrop ND-1000 spectrophotometer at 260 nanometers (Thermo Fisher Technologies Inc.,

499    Santa Clara, CA, USA), diluted to one nanogram/microliter, and ten nanograms used to prepare

500    sequencing libraries as described below.  Sequence data for the three patient samples were

501    deposited to the NCBI Sequence Read Archive in project accession PRJNA521396.

502

503    **DNA library preparation and sequencing**

504        DNA libraries were prepared and sequenced for all samples utilizing Ion Torrent reagents

505    and the Ion Torrent Proton sequencer (Thermo Fisher Technologies Inc., Santa Clara, CA, USA).

506    Ten nanograms of DNA was input to the Ion Xpress Plus Fragment Library Kit (manual

507    #MAN0009847, revC). DNA was sheared using the Ion Shear enzymatic reaction for 12 min,

508    and Ion Xpress barcode adapters were ligated following end repair. Resulting libraries were

509    amplified using the manufacturer supplied library amplification primers and recommended

510    conditions. Amplified libraries were size selected to approximately 200 base pairs using E-gel

511    SizeSelect Agarose cassettes (Invitrogen, Carlsbad, CA, USA) as outlined in the Ion Xpress

512    manual and quantitated with the Ion Universal Library quantitation kit. Equimolar amounts of

513    the library were templated with an Ion PI Template OT2 200 kit V3. The resulting templated

514    beads were enriched with the Ion OneTouch ES system and quantitated with the Qubit Ion

515    Sphere Quality Control kit on a Qubit 3.0 fluorimeter (Qubit, NY, NY, USA). Enriched

516    templated beads were loaded onto an Ion PI V2 chip and sequenced according to the

517    manufacturer's protocol using the Ion PI Sequencing 200 kit V3. Data were processed with Ion

518    Torrent Server software v4.4.3 to produce data files in BAM format.

519

520    **Read processing and quality control**

521         Sequences were converted to FASTQ format from raw BAM files with bedtools'

522    bamtofastq (Quinlan & Hall, 2010)2.17.0, (Quinlan & Hall, 2010). FastQC ("Babraham

523    Bioinformatics - FastQC A Quality Control tool for High Throughput Sequence Data," n.d.)

524    v0.11.5, ("Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput

525    Sequence Data," n.d.) was used to generate sequence quality reports. FastX toolkit (Gordon &

526    Hannon, 2010)v.0.0.14, (Gordon & Hannon, 2010) was used to perform quality control measures

527    on FASTQ data including quality filtering, trimming, setting a minimum read length, and

528    removal of duplicate reads. Files were converted to FASTA with FastX. Data files before and

529    after QC were used as input to Centrifuge when testing the effect of quality control; otherwise,

530    all files were quality controlled before analysis.

531

532    **Removing host contamination by aligning to the human genome**

533        To remove host (human) reads, FASTQ read files were mapped to HG38 (Genome

534    reference consortium human genome build38) using Bowtie2 (Langmead & Salzberg, 2012)

535    using the --very-sensitive option. Human reads were removed by alignment from patient data

536    before the analysis in Centrifuge except when testing the effect of host screening by other

537    methods.

538

539    **Centrifuge and CLARK read classification**

540        CLARK v1.1.3 (Ounit et al., 2015a) was used to classify reads to known taxa using the

541    default CLARK database and parameters. Centrifuge v1.0.3-beta (Kim et al., 2016) was used to

542    classify reads to known taxa with a custom database generated from 23,276 complete archaeal,

543    bacterial, and viral genomes downloaded from Refseq in July 2017 using the centrifuge-

544    download and centrifuge-build scripts respectively. The custom database is available at

545    https://github.com/hurwitzlab/NeutropenicFever.

546

547    **Binary mixture Centrifuge results filtering**

548        Centrifuge abundance report results were filtered to only include organisms at the species

549    or strain-level with a minimum of 0.1% of total reads classified and at least 0.05% abundance as

550    calculated by Centrifuge. These settings were chosen based on the known abundances used in the

551    mixtures. False positive was calculated by summing the relative abundances of any organism

552    identified by Centrifuge that was not added to the mixture. Centrifuge reports read-matches to

553    phage separately from their host species; however, no phage or prophage passed the above

554    filters, so there was no effect on the relative abundance calculations for the binary mixtures. The

30

555     coefficient of determination ($R^2$) was calculated based on the log of both relative abundance

556     estimates at each known dilution.

557

558     **Mock community Centrifuge results filtering**

559     Centrifuge abundance report results were filtered to only include organisms at the

560     species or strain-level with a minimum of at least 0.005% abundance as calculated by Centrifuge

561     and no minimum number of reads. These settings were chosen based on the known abundances

562     calculated for the mock community which was lower than the bacterial mixtures (0.01%). In the

563     case of the mock community, two species-specific phages were identified that passed the filters

564     (*Pseudomonas* phage with relative abundance 1.5%, and *Staphylococcus* phage with relative

565     abundance 0.8%). The matches to these phages were included when calculating relative

566     abundances for the 20 organisms, but not included in the figure. The coefficient of determination

567     ($R^2$) was calculated based on the log of the relative abundance estimates for all 20 species.

568

569

570     **Febrile Neutropenia Centrifuge results filtering**

571     Centrifuge abundance report results were filtered to only include organisms at the species

572     or strain-level with a minimum of 1% of total reads classified and at least 5% abundance as

573     calculated by Centrifuge. Similarly to the bacterial mixtures, no phage or prophage passed the

574     filters above, so there was no effect on relative abundance calculations.

575

576     **Genome coverage of suspected pathogens from febrile neutropenia patient samples**

577    To determine genome coverage, we used Bowtie2 (Langmead & Salzberg, 2012) to map

578    FASTQ reads (with option --very-sensitive) to reference genomes for the organisms identified by

579    Centrifuge (*Pseudomonas fluorescens* accession: NC_012660.1, Human Parvovirus B19

580    accession: NC_000883.2, Torque Teno Virus accession: NC_015783.1). Resulting BAM files

581    were then analyzed utilizing Samtools' (v1.3.1, (Li et al., 2009) depth tool to generate coverage

582    values and visualized in R v3.1.1 (R scripts are available here:

583    https://github.com/hurwitzlab/NeutropenicFever ).

584

585    **Software availability**

586    To improve access to Centrifuge and the bubble chart visualizations used in this

587    manuscript, both tools have been made available on iMicrobe (https://www.imicrobe.us). As a

588    starting point, researchers may run centrifuge-0.0.6u1 followed by centrifuge-bubble-0.0.5u1 to

589    reproduce the bacterial mixing results in the manuscript using the sample data provided. Source

590    code for running centrifuge on a high-performance compute cluster is available in Github at

591    https://github.com/hurwitzlab/Centrifuge_HPC and analyses, scripts and visualizations are also

592    archived at https://github.com/hurwitzlab/NeutropenicFever.

593

594    **ACKNOWLEDGMENTS**

600

608
609  **REFERENCES**

610  Ashton, P. M., Peters, T., Ameh, L., McAleer, R., Petrie, S., Nair, S., … Dallman, T. (2015). Whole

611    Genome Sequencing for the Retrospective Investigation of an Outbreak of Salmonella Typhimurium

612    DT 8. *PLoS Currents*, *7*.

613    https://doi.org/10.1371/currents.outbreaks.2c05a47d292f376afc5a6fcdd8a7a3b6

614  Babraham Bioinformatics - FastQC A Quality Control tool for High Throughput Sequence Data. (n.d.).

615    Retrieved May 23, 2018, from https://www.bioinformatics.babraham.ac.uk/projects/fastqc/

616  Bazinet, A. L., & Cummings, M. P. (2012). A comparative evaluation of sequence classification

617    programs. *BMC Bioinformatics*, *13*, 92.

618  de Villiers, E.-M., Borkosky, S. S., Kimmel, R., Gunst, K., & Fei, J.-W. (2011). The diversity of torque

619    teno viruses: in vitro replication leads to the formation of additional replication-competent subviral

620    molecules. *Journal of Virology*, *85*(14), 7284–7295.

621  Ecker, D. J., Sampath, R., Li, H., Massire, C., Matthews, H. E., Toleno, D., … Tang, Y.-W. (2010). New

622    technology for rapid molecular diagnosis of bloodstream infections. *Expert Review of Molecular*

623    *Diagnostics*, *10*(4), 399–415.

624  Frey, K. G., Herrera-Galeano, J. E., Redden, C. L., Luu, T. V., Servetas, S. L., Mateczun, A. J., …

625    Bishop-Lilly, K. A. (2014). Comparison of three next-generation sequencing platforms for

626    metagenomic sequencing and identification of pathogens in blood. *BMC Genomics*, *15*, 96.

627    Gordon, A., & Hannon, G. J. (2010). Fastx-toolkit. *FASTQ/A Short-Reads Preprocessing Tools*

628    *(unpublished) Http://hannonlab. Cshl. Edu/fastx_toolkit*.

629    Kim, D., Song, L., Breitwieser, F. P., & Salzberg, S. L. (2016). Centrifuge: rapid and sensitive

630    classification of metagenomic sequences. *Genome Research*, *26*(12), 1721–1729.

631    Land, M., Hauser, L., Jun, S.-R., Nookaew, I., Leuze, M. R., Ahn, T.-H., … Ussery, D. W. (2015).

632    Insights from 20 years of bacterial genome sequencing. *Functional & Integrative Genomics*, *15*(2),

633    141–161.

634    Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*,

635    *9*(4), 357–359.

636    Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform.

637    *Bioinformatics* , *25*(14), 1754–1760.

638    Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., … 1000 Genome Project Data

639    Processing Subgroup. (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics* ,

640    *25*(16), 2078–2079.

641    Liu, H., Fu, Y., Xie, J., Cheng, J., Ghabrial, S. A., Li, G., … Jiang, D. (2011). Widespread endogenization

642    of densoviruses and parvoviruses in animal and human genomes. *Journal of Virology*, *85*(19), 9863–

643    9876.

644    M. Burrows, D. J. W. (1994). A block-sorting lossless data compression algorithm. Retrieved from

645    http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.3.8069

646    McIntyre, A. B. R., Ounit, R., Afshinnekoo, E., Prill, R. J., Hénaff, E., Alexander, N., … Mason, C. E.

647    (2017). Comprehensive benchmarking and ensemble approaches for metagenomic classifiers.

648    *Genome Biology*, *18*(1), 182.

649    Mollerup, S., Friis-Nielsen, J., Vinner, L., Hansen, T. A., Richter, S. R., Fridholm, H., … Hansen, A. J.

650    (2016). Propionibacterium acnes: Disease-Causing Agent or Common Contaminant? Detection in

34

651     Diverse Patient Samples by Next-Generation Sequencing. *Journal of Clinical Microbiology*, *54*(4),

652     980–987.

653 Naccache, S. N., Federman, S., Veeraraghavan, N., Zaharia, M., Lee, D., Samayoa, E., … Chiu, C. Y.

654     (2014). A cloud-compatible bioinformatics pipeline for ultrarapid pathogen identification from next-

655     generation sequencing of clinical samples. *Genome Research*, *24*(7), 1180–1192.

656 Ounit, R., Wanamaker, S., Close, T. J., & Lonardi, S. (2015a). CLARK: fast and accurate classification of

657     metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics*, *16*, 236.

658 Ounit, R., Wanamaker, S., Close, T. J., & Lonardi, S. (2015b). CLARK: fast and accurate classification of

659     metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics*, *16*, 236.

660 Parize, P., Muth, E., Richaud, C., Gratigny, M., Pilmis, B., Lamamy, A., … Eloit, M. (2017). Untargeted

661     next-generation sequencing-based first-line diagnosis of infection in immunocompromised adults: a

662     multicentre, blinded, prospective study. *Clinical Microbiology and Infection: The Official*

663     *Publication of the European Society of Clinical Microbiology and Infectious Diseases*.

664     https://doi.org/10.1016/j.cmi.2017.02.006

665 Park, H. J., Na, S., Park, S. Y., Moon, S. M., Cho, O.-H., Park, K.-H., … Choi, S.-H. (2011). Clinical

666     significance of Propionibacterium acnes recovered from blood cultures: analysis of 524 episodes.

667     *Journal of Clinical Microbiology*, *49*(4), 1598–1601.

668 Quick, J., Loman, N. J., Duraffour, S., Simpson, J. T., Severi, E., Cowley, L., … Carroll, M. W. (2016).

669     Real-time, portable genome sequencing for Ebola surveillance. *Nature*, *530*(7589), 228–232.

670 Quinlan, A. R., & Hall, I. M. (2010). BEDTools: a flexible suite of utilities for comparing genomic

671     features. *Bioinformatics* , *26*(6), 841–842.

672 Rosen, G. L., Reichenberger, E. R., & Rosenfeld, A. M. (2011). NBC: the Naive Bayes Classification tool

673     webserver for taxonomic classification of metagenomic reads. *Bioinformatics* , *27*(1), 127–129.

674 Schmieder, R., & Edwards, R. (2011). Fast identification and removal of sequence contamination from

675     genomic and metagenomic datasets. *PloS One*, *6*(3), e17288.

676 Snitkin, E. S., Zelazny, A. M., Thomas, P. J., Stock, F., NISC Comparative Sequencing Program Group,

677     Henderson, D. K., … Segre, J. A. (2012). Tracking a hospital outbreak of carbapenem-resistant

678     Klebsiella pneumoniae with whole-genome sequencing. *Science Translational Medicine*, *4*(148),

679     148ra116.

680  van Walraven, C., & Wong, J. (2014). Independent influence of negative blood cultures and bloodstream

681     infections on in-hospital mortality. *BMC Infectious Diseases*, *14*, 36.

682  Watts, G. S., Youens-Clark, K., Slepian, M. J., Wolk, D. M., Oshiro, M. M., Metzger, G. S., … Hurwitz,

683     B. L. (2017). 16S rRNA gene sequencing on a benchtop sequencer: accuracy for identification of

684     clinically important bacteria. *Journal of Applied Microbiology*, *123*(6), 1584–1596.

685  Wilson, M. R., Naccache, S. N., Samayoa, E., Biagtan, M., Bashir, H., Yu, G., … Chiu, C. Y. (2014).

686     Actionable diagnosis of neuroleptospirosis by next-generation sequencing. *The New England*

687     *Journal of Medicine*, *370*(25), 2408–2417.

688  Wohlfarth, P., Leiner, M., Schoergenhofer, C., Hopfinger, G., Goerzer, I., Puchhammer-Stoeckl, E., &

689     Rabitsch, W. (2018). Torquetenovirus Dynamics and Immune Marker Properties in Patients

690     Following Allogeneic Hematopoietic Stem Cell Transplantation: A Prospective Longitudinal Study.

691     *Biology of Blood and Marrow Transplantation: Journal of the American Society for Blood and*

692     *Marrow Transplantation*, *24*(1), 194–199.

693  Wong, V., Levi, K., Baddal, B., Turton, J., Boswell, T.C. (2011). Spread of *Pseudomonas fluorscens* Due

694     to Contaminated Drinking Water in a Bone Marrow Transplant Unit. *Journal of Clinical*

695     *Microbiology*, *49(6)*, 2093-2096.

696  Wood, D. E., & Salzberg, S. L. (2014). Kraken: ultrafast metagenomic sequence classification using exact

697     alignments. *Genome Biology*, *15*(3), R46.

698
699
700
701
702
703
704

705

706

707

708

709

710

711

712

**SUPPORTING INFORMATION**

714

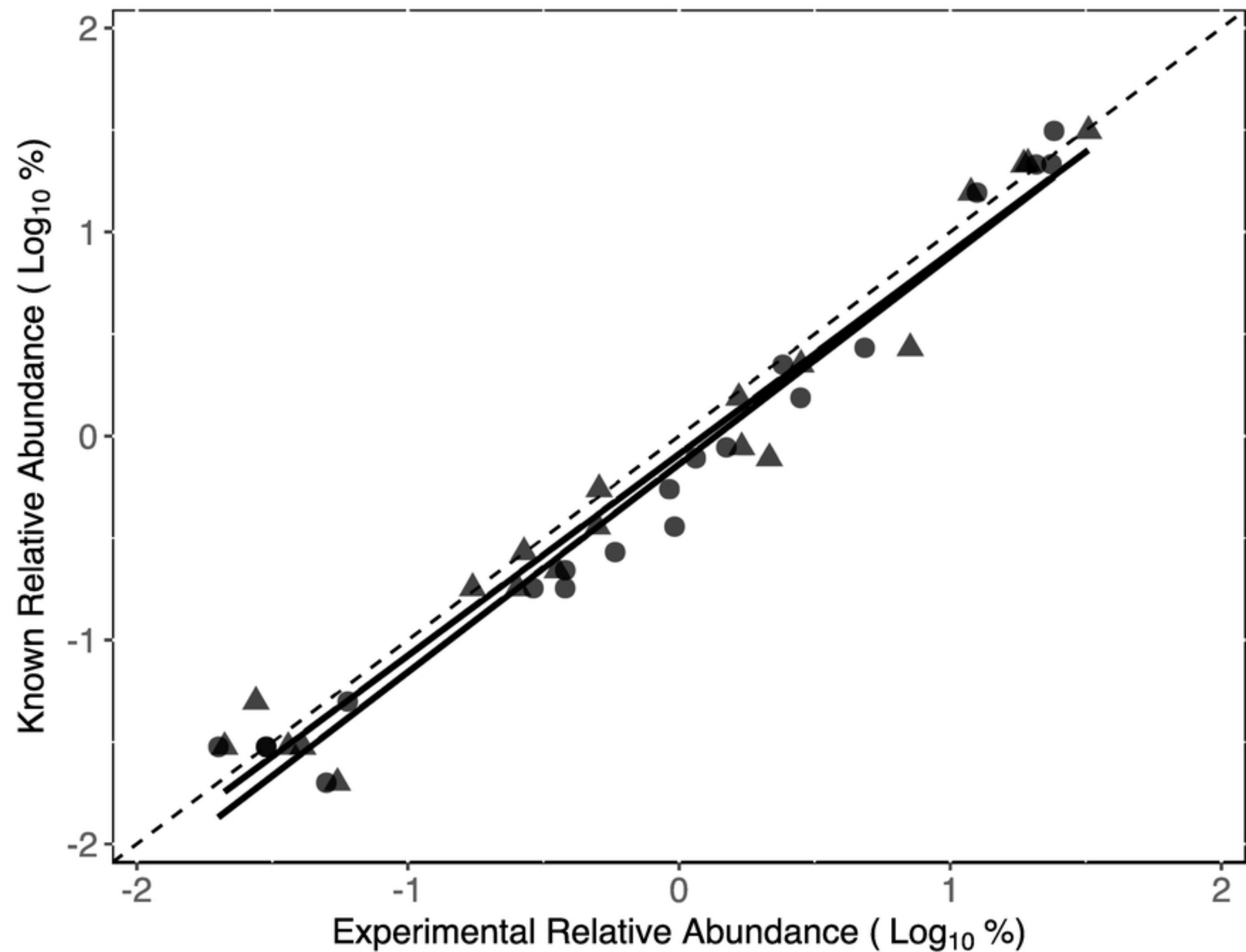**S1. Quality Control Reduction of Reads**
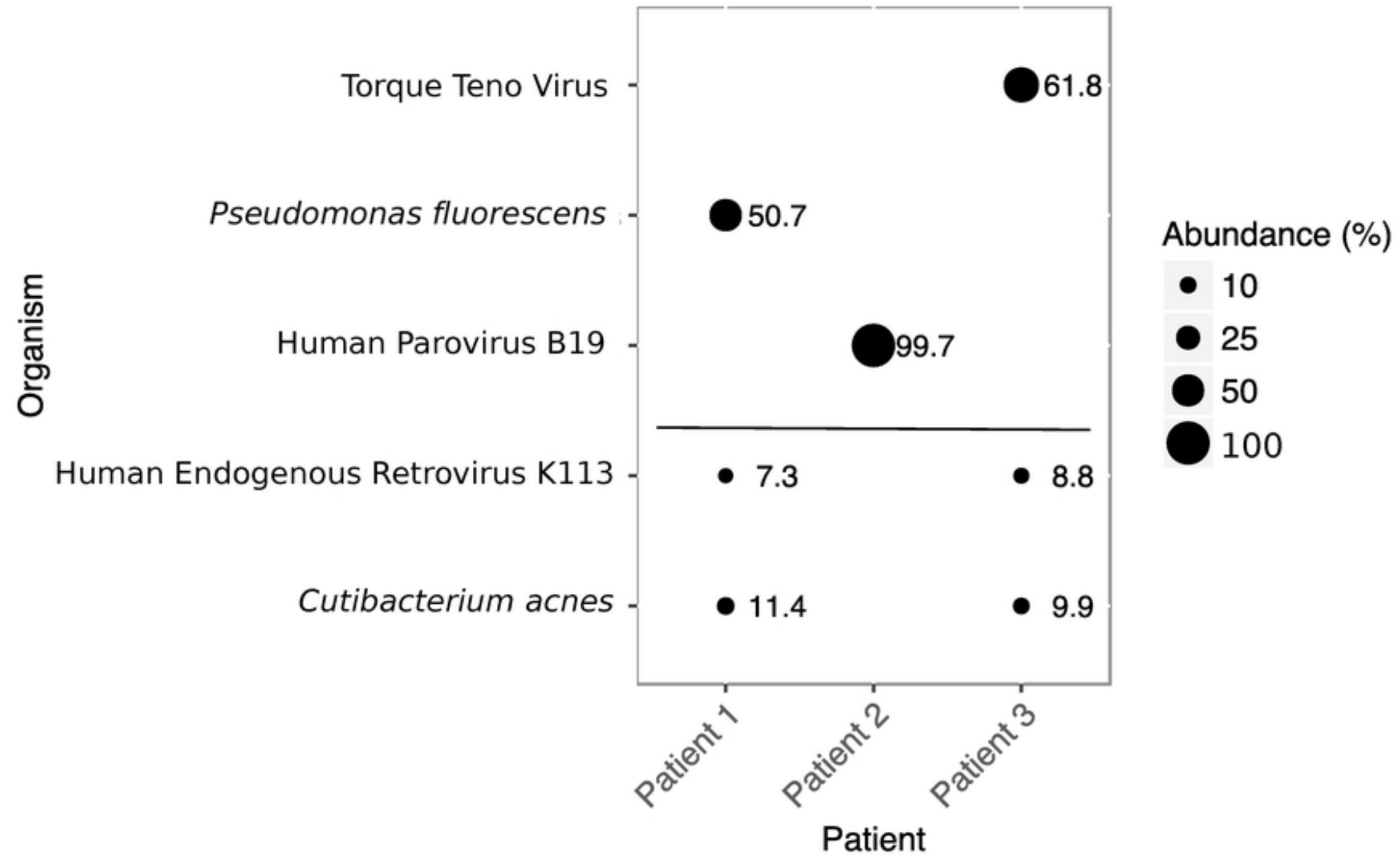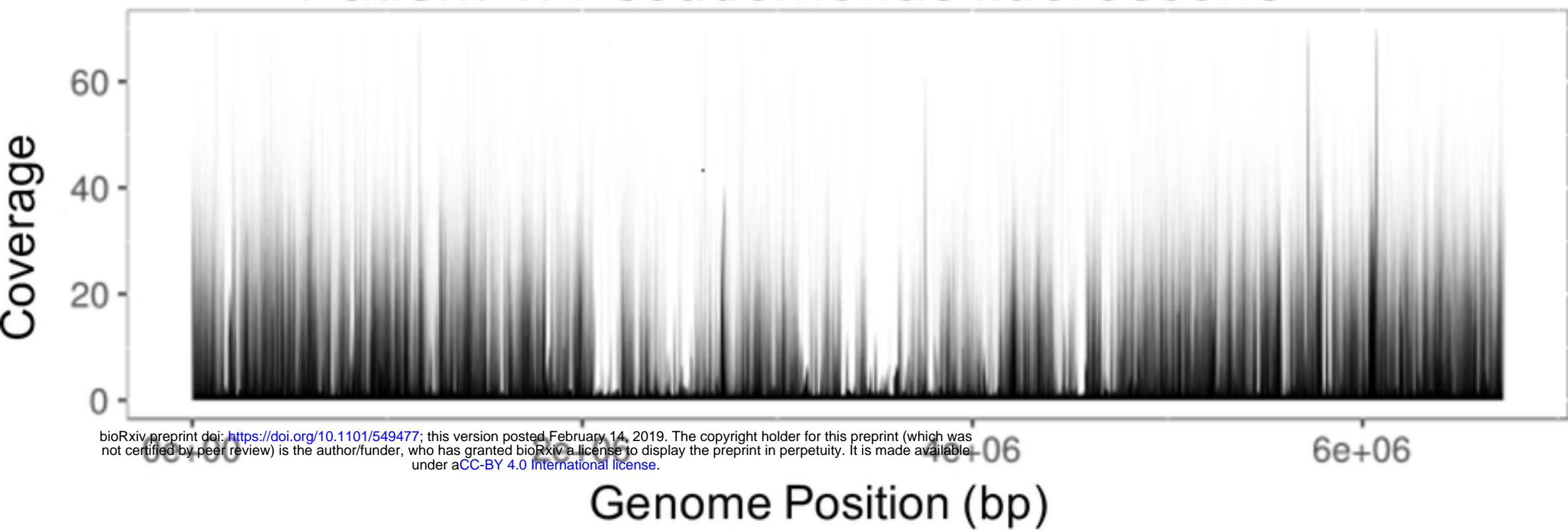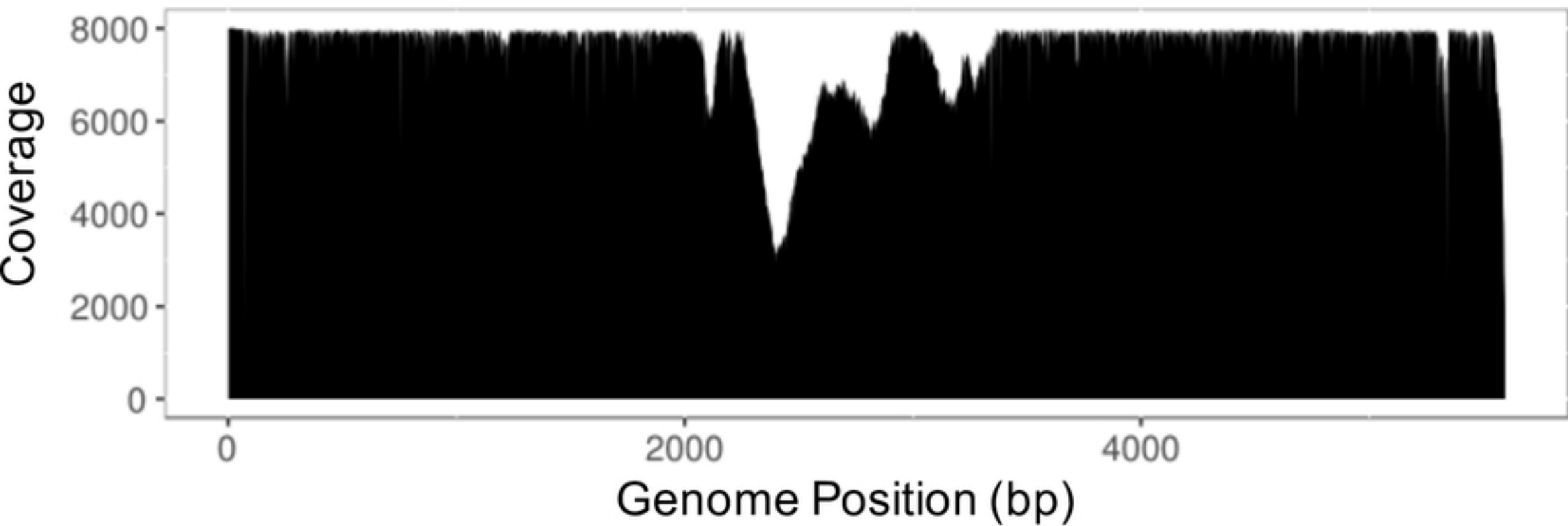
Figure 1

Figure 2

Figure 3

# Patient 1: *Pseudomonas fluorescens*

# Patient 2: Human Parvovirus



# Patient 3: Torque Teno Virus



Figure 4

Figure 5

Figure 6