1 **Title:** "Evaluating the Benefits and Limits of Multiple Displacement Amplification with Whole-
2 Genome Oxford Nanopore Sequencing"
3
4 **Running Title:  Evaluating MDA-ONT Genome Sequencing**
5

6 **Authors:**
7 Fiifi A. Dadzie [1], Megan S. Beaudry [2*], Alex Deyanov [3], Haley Slanis [3], Minh Q. Duong [3], Randi
8 Turner [4,5], Asis Khan [5], Cesar A. Arias [3,6,7], Jessica C. Kissinger [1,4,8], Travis C. Glenn [1,2,8], Rodrigo de
9 Paula Baptista [3,4,6,7,8]
10
11 **Addresses:**
12 [1]Department of Genetics, University of Georgia, Athens, GA USA 30602;
13 [2] Department of Environmental Health Science, University of Georgia, Athens, GA USA 30602
14 [3] Center for Infectious Disease, Houston Methodist Research Institute, Houston, TX USA 77030;
15 [4] Center for Tropical and Emerging Global Diseases, University of Georgia, Athens, GA USA;
16 [5] USA Department of Agriculture,  Agricultural Research Service, Beltsville Agricultural Research
17 Service, Animal Parasitic Disease Laboratory, Beltsville, MD USA;
18 [6] Division of Infectious Diseases and Department of Medicine, Houston Methodist Hospital,
19 Houston, TX USA 77030;
20 [7]Department of Medicine, Weill Cornell Medical College, New  York, NY
21 [8] Institute of Bioinformatics, University of Georgia, Athens, GA USA 30602;
22
23 [*]Current address: Daicel Arbor Biosciences, 5840 Interface Dr. Suite 101, Ann Arbor, MI 48103;
24
25
26 Corresponding author:  Rodrigo de Paula Baptista
27                        Houston Methodist Research Institute, Houston, TX USA 77030
28                        rdepaulabaptista@houstonmethodist.org
29
30

**ABSTRACT**

Multiple Displacement Amplification (MDA) outperforms conventional PCR in long fragment and whole genome amplification which makes it attractive to couple with long-read sequencing of samples with limited quantities of DNA to obtain improved genome assemblies. Here, we explore the efficacy and limits of MDA for genome sequence assembly using Oxford Nanopore Technologies (ONT) rapid library preparations and minION sequencing. We successfully generated almost complete genome sequences for all organisms examined, including *Cryptosporidium meleagridis*, *Staphylococcus aureus*, *Enterococcus faecium*, and *Escherichia coli*, with the ability to generate high-quality data from samples starting with only 0.025 ng of total DNA. Controlled sheared DNA samples exhibited a distinct pattern of size-increase after MDA, which may be associated with the amplification of long, low-abundance fragments present in the assay, as well as generating concatemeric sequences during amplification. To address concatemers, we developed a computational pipeline (CADECT: Concatemer Detection Tool) to identify and remove putative concatemeric sequences. This study highlights the efficacy of MDA in generating high-quality genome assemblies from limited amounts of input DNA. Also, the CADECT pipeline effectively mitigated the impact of concatemeric sequences, enabling the assembly of contiguous sequences even in cases where the input genomic DNA was degraded. These results have significant implications for the study of organisms that are challenging to culture *in vitro*, such as *Cryptosporidium*, and for expediting critical results in clinical settings with limited quantities of available genomic DNA.

**Key Words:** Infectious Diseases, Apicomplexa, Enterococcaceae, Enterobacterales, low abundance DNA, LRS Special Issue.

**INTRODUCTION**

The advent of next-generation sequencing technologies has revolutionized genomics research by enabling the rapid and cost-effective generation of vast amounts of sequencing data (Slatko et al. 2018; Hu et al. 2021). Among these technologies, Oxford Nanopore Sequencing (ONT) stands out due to its ability to provide long-read sequencing data in real-time, with lower instrument costs and less input DNA required for non-amplified library preparations than the other major commercial long-read sequencing platform, PacBio (Pacbio 2022). ONT sequencing has been used for numerous applications, including *de novo* genome assembly, metagenomics, and pathogen detection. However, ONT sequencing library preparations typically still requires higher-quality and higher quantities of DNA inputs than may be available for many projects. ONT rapid

67    library preparations usually require at least 50 ng input per sample, but more is required when
68    pooling with <8 other barcoded samples (≥400 ng is recommended for loading onto a MinION
69    flow cell). Many samples also suffer from DNA degradation, where the majority of DNA fragments
70    are shorter than is desirable for ONT library preparation and removal of small fragments further
71    reduces the quantity of DNA available. This poses challenges when working with samples that
72    have limited quantity and/or degraded DNA (Delahaye and Nicolas 2021). For this reason,
73    alternative library preparation or sequencing techniques, including short-read sequencers (*e.g.*,
74    Illumina, Element Biosciences AVITI), are often preferred for handling samples with low
75    molecular weight and/or low quantities of DNA.
76    To overcome these limitations, multiple displacement amplification (MDA) has emerged as a
77    valuable and highly efficient method for amplifying small quantities of DNA. MDA has significant
78    advantages over conventional PCR and other whole genome amplification techniques (Hou et al.
79    2015). These advantages include reduced waste of rare samples, isothermal amplification for
80    efficiency, heightened sensitivity in detecting low amounts of DNA inputs, minimized bias and
81    error rates, amplification of long DNA fragments and whole genome amplification of organisms
82    with relatively small genome size (< 10Mb). MDA utilizes the Phi29 DNA polymerase with a
83    displacement activity that enables the isothermal amplification of DNA with high fidelity and
84    exponential amplification of DNA molecules (Dean et al. 2002). This technique has been
85    successfully applied in various genomic studies, including single-cell sequencing, ancient DNA
86    analysis, and microbiome studies (Binga et al. 2008; Lasken 2009). Moreover, MDA enables the
87    amplification of long DNA fragments, making it valuable for applications such as cloning and
88    genomic library preparation (Fullwood et al. 2008). While a protocol for MDA with ligation
89    sequencing kits (Qiagen, Germany) is available, MDA's application with ONT rapid kits, which
90    offer faster processing times and yield relatively smaller fragments compared to ligation kits, has
91    not been extensively investigated. Consequently, MDA's potential limitations and impacts on
92    whole-genome assembly in this context remain relatively unexplored.
93    The use of MDA combined with ONT sequencing has the potential to unlock genomic insights for
94    organisms that are small (e.g., larval ticks, parasitoid wasps, etc.) to microscopic, especially those
95    that are difficult or impossible to culture *in vitro* (*e.g*. *Cryptosporidium* species, *Mycobacterium*
96    *leprae* and *Treponema pallidum*). Furthermore, clinical samples and isolates with limiting
97    amounts of DNA pose a challenge for rapid and accurate genome sequence analysis, especially
98    in urgent clinical situations where timely results are crucial. Working with degraded DNA samples
99    becomes an issue since it could limit the sequence genomic coverage and assembly (Ceccherini
100   et al. 2003). MDA is not suitable for analysis of severely degraded DNA, since could impact: (i)
101   MDA efficiency due potential breaks or lesions leading incomplete or suboptimal amplification;
102   (ii) bias resulting in uneven coverage across the genome; and (iii) contaminants that could
103   interfere with the MDA reaction (Wang et al. 2004).

104  It's important to mention that when utilizing MDA there are limitations that needs to be
105  considered to ensure the reliability and integrity of the sequencing results. While MDA has
106  facilitated genomic sequencing from low concentrations of template nucleic acid, there are still
107  several limitations to consider. These include: (i) Nonspecific amplification resulted from primer
108  dimer formation causing template switching or contamination by DNA templates; (ii) Formation
109  of chimeric DNA rearrangements; and (iii) Representation bias, which can affect the accuracy and
110  completeness of the amplified genomic material (Binga et al. 2008). Some studies shows that
111  chimeric reads are usually invert chimeras or direct chimeras, but it was previously observed that
112  most of detected MDA chimeric sequences (85%) are inverted chimeras, such as inverted
113  sequences with intervening deletions which can be caused by template switching (Lasken and
114  Stockwell 2007; Lu et al. 2023). These chimeric sequences are known to affect genome
115  sequencing since they can be considered as amplification artifacts, which cannot be used for
116  genome assembling (Lu et al. 2023). Studies suggest that chimerism in MDA sequencing data is a
117  significant concern that is gaining attention, particularly with the rise of single-cell studies (Hard
118  et al. 2023).
119  To address the challenges associated with artifactual concatemeric sequences generated during
120  MDA, we developed a novel bioinformatic tool called CADECT (Concatemer Detection Tool),
121  which is made available at https://github.com/rpbap/CADECT. This tool enabled the
122  identification and removal of putative inverted chimeric concatemers, thus improving the
123  accuracy and contiguity of the genome assembly.
124  Our study aims to provide valuable insights into the use of MDA for whole-genome ONT
125  sequencing, particularly for low molecular weight and/or low quantities of DNA samples,
126  highlighting its potential as a powerful method to obtain high-quality long-read sequencing data.
127  We assessed the MDA advantages and constraints, and effectiveness for whole-genome
128  assembly in microbial organisms with genome sizes <10 Mb. This is especially significant for
129  infectious disease agents, where obtaining enough DNA can be challenging. Overall, our study
130  underscores the potential of MDA in enabling high-quality long-read sequencing from challenging
131  low-concentration DNA samples, emphasizing its importance in various genomic research and
132  clinical applications.
133

## RESULTS

135

**WGA** results

137    Our whole genome amplification (WGA) results reveal that in each sample type tested, we find

138 an overall fold change of > 500× in comparison to the original sample (Table 1). Following

139 amplification, approximately 1.5 µg of the product was debranched using T7 endonuclease

140 prior to library preparation for ONT sequencing. Typically, we experience a ~45% recovery after

141 this step, attributed to the bead purification process (Table 1). Though a significant amount of

142 DNA is lost during the DNA purification step post T7 endonuclease reaction, an overall fold

143 change of ~100× is observed when compared to the WGA DNA input.

144

145 **Table 1 - Observed amplification yield increase by sample type**

| | WGA input (ng) | WGA output (ng) | T7 output (ng) | T7 recovery (%) | Estimated Fold Increase |
|---|---|---|---|---|---|
| Gram-positive (*S. aureus*) | 2.5 | 1500 | 894 | 59.6 | 357.6× |
| Gram-negative (*E. coli*) | 5.0 | 8360 | 552 | 36.8 | 110.0× |
| Eukaryotic Pathogen (*Cryptosporidium* ssp.) | 2.5 | 1976 | 555 | 44.1 | 222.0× |
| Background (Calf thymus) | 5.0 | 3800 | 620 | 41.33 | 124.0× |

146

147 We successfully obtained contiguous and sometimes even chromosomal-level assemblies from

148 the samples analyzed in this study, starting with DNA inputs significantly lower than Oxford

149 Nanopore's recommended minimum of 50 ng for the rapid barcode kit (Table S1).

150

151 For certain samples, such as *E. faecium*, we observed that achieving improved contiguity required

152 generating higher depth coverage during the sequencing. Our results indicate that, for this

153 organism, reaching depths beyond 70× allowed us to attain a chromosomal-level assembly with

154 only 2.5 ng of starting total DNA (Table S2). In comparison, increased sequencing depth on

155 samples that started with less than 0.001 ng of input into the MDA did not enhance contiguity.

156 Combining separate MDA amplifications of the same limited input samples did improve the final

157 genome coverage because the random nature of the initial templates and amplification process.

158 Thus, multiple independent MDAs appears to be advantageous because it could randomly

159 amplify by chance different regions that are beneficial for the genome assembly.

160 To check for potential GC bias on the sequencing depth along the genome, the $R^2$ correlation

161 coefficient between average depth and average %GC across 1000 base pair regions of the *E.*

162 *faecium* non-amplified and amplified assemblies was 0.0262 and 0.0265, respectively (Fig. S1).

5

163

164 **WGA amplification using serially diluted samples**

165

166 Serial dilutions of a single *E. faecium* sample reveal successful DNA amplification even at low
167 initial DNA amount of 2.5E-5 ng (Fig. 1A). The MDA technique imposes a size limit on its amplified
168 products, with an average product length of 10-12 kb (Dean et al. 2002), and it requires a
169 debranching step, leading to a reduction in the mean sequence read sizes (Table S3). Post MDA,
170 the average size of the reads typically falls within the range of 2-3 kb. In contrast, standard Oxford
171 Nanopore Technologies (ONT) assays without amplification (*i.e.*, ONT Rapid Barcode Kit (RBK))
172 which includes a transposase step that simultaneously cleaves template molecules and attaches
173 tags to the cleaved ends, typically generate DNA fragments ranging from 5-20 kb. However, when
174 assessing genome coverage (genome sizes <10 Mb), we observed that DNA inputs below 0.025
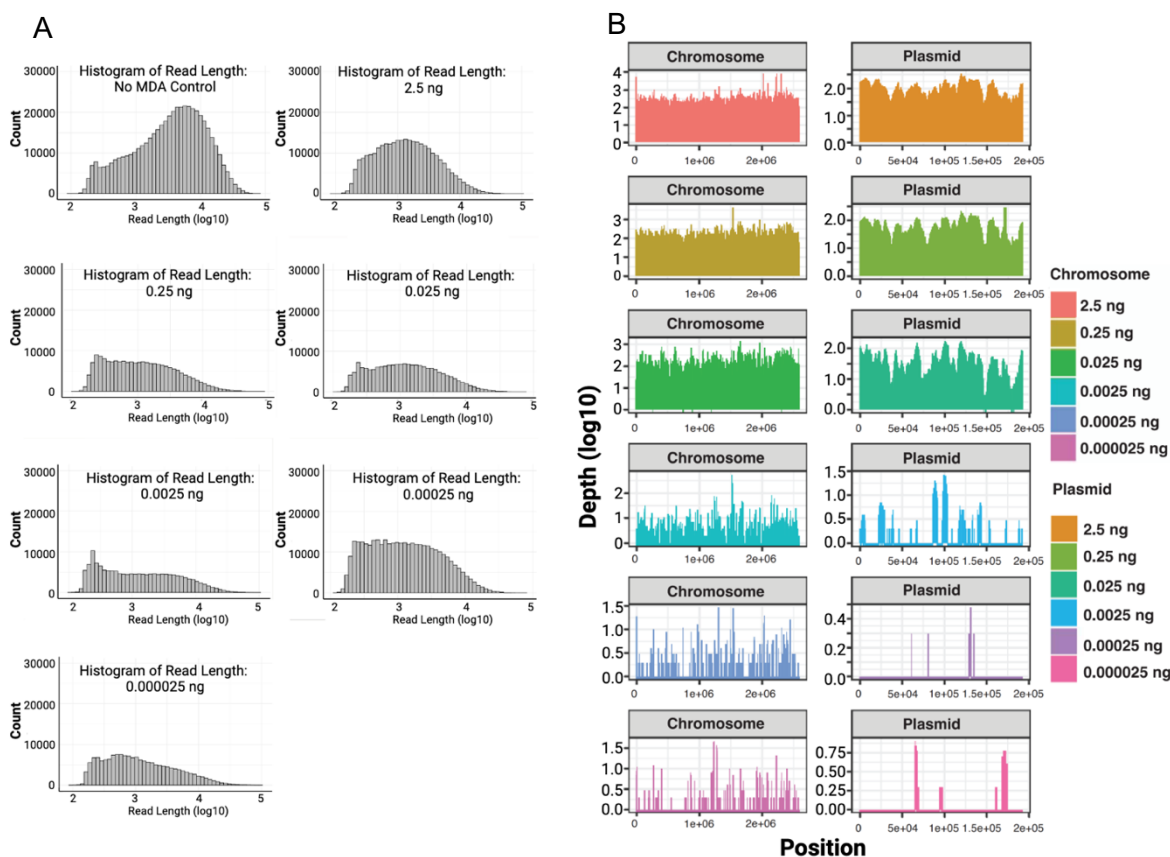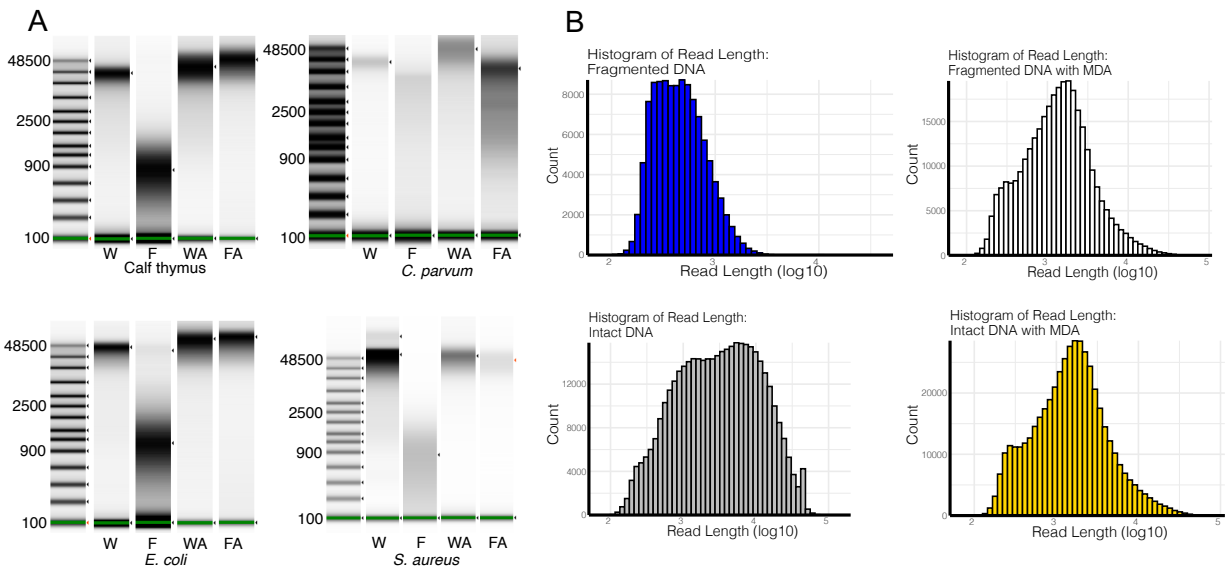175 ng result in incomplete coverage of certain genomic regions (Fig. 1B).

176



196 **Figure 1**. **Serial dilution test with *E. faecium* sample DNA**. (A) Depicts the distribution and
197 counts (y-axis) of read lengths (x-axis) across all diluted samples. (B) Illustrates the horizontal
198 coverage of the chromosomal regions across all diluted and subsequently amplified samples,
199 with read depth on the y-axis and genome position on the x-axis.

200

**MDA results in an unexpected size increase from fragmented DNA samples**

Following controlled enzymatic fragmentation using a dsDNA fragmentase and MDA according to our protocol, we observed an unexpected size increase distribution of fragments (Fig. 2A). Indeed, for all samples, except *Cryptosporidium*, the size distribution post-MDA was nearly identical for intact and fragmented input DNA. Subsequent analysis of the ONT sequencing results revealed the existence of read lengths that are longer than those present in the same sample without MDA amplification (Fig. 2B).



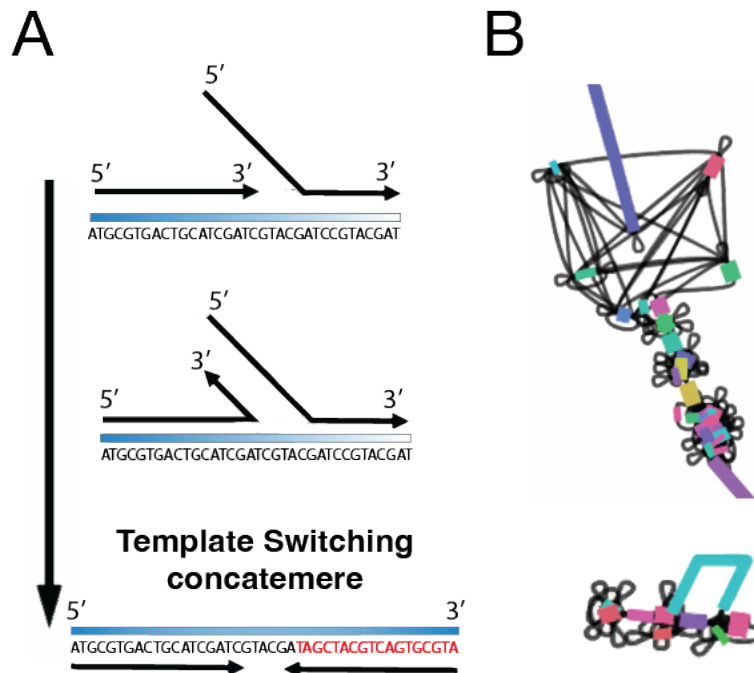**Figure 2. DNA fragment and read Size Range pre- and post- Multiple Displacement Amplification using Size-Controlled Fragmented DNA.** (A) TapeStation results for different organism sets; (B) ONT sequencing results obtained before and after amplification for *S. aureus*. W = whole intact DNA; F = Fragmented DNA; WA and FA = after amplification; and WT and FT = after T7 debranching. Uncropped TapeStation results are in Fig. S2.

Upon closer examination of the sequence content, two distinct types of reads were identified. Some represented potentially low-abundance longer reads that escaped fragmentation during the enzyme incubation and were subsequently amplified. The other reads were primarily chimeric concatemers, likely generated through template switching of short fragments during MDA (Fig. 3A). While the occurrence of concatemers in MDA assays has been reported previously (Paul and Apgar 2005; Lu et al. 2023), they are typically present in low amounts after sequencing. In our case, the fragmentation process seemed to enhance the prevalence of these chimeric reads in our ONT sequencing. As expected, assembly of the data revealed that the presence of

227 these chimeric/concatemers regions significantly impacted genome assembly, resulting in bubble
228 fragmentation effect across the entire genome and affecting contiguity (Fig. 3B).
229



230
231

**Figure 3. Impact of MDA-Generated Concatemers on the Genome Assembly**. (A) Concatemers generated by template switching; (B) Graph representation of the effect of concatemers on genome assembly (bubble fragmentation effect)
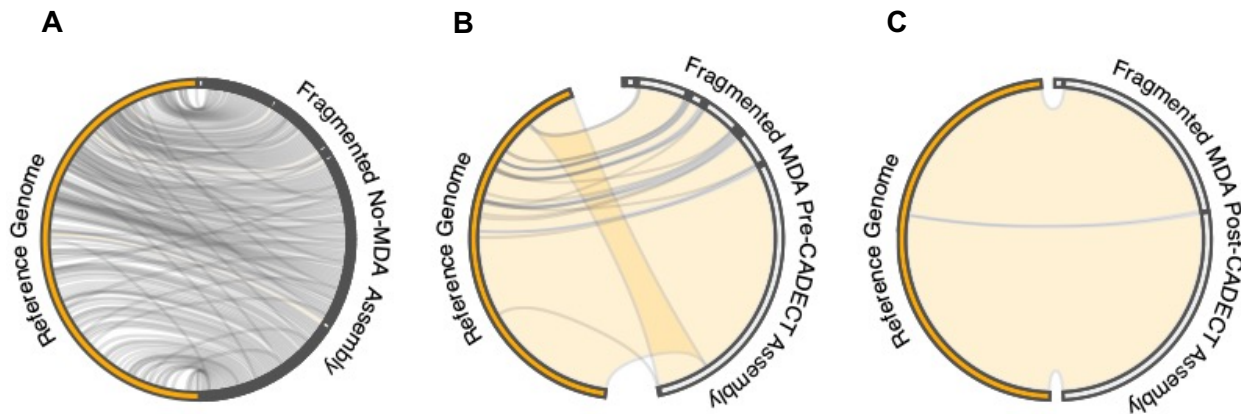
**Concatemer detection tool**

To identify and eliminate potential concatemers generated by MDA, we designed a concatemer detection tool specifically tailored for raw ONT reads called CADECT. This tool enables the differentiation of putative concatemeric chimeric reads from non-concatemeric ones. To achieve this, the process involves dividing each long-read sequence into multiple fragments using a sliding window approach and then aligning these fragments with one another. The underlying hypothesis is that the presence of a concatemer would result in certain windows aligning with each other, thereby confirming the existence of a potential concatemer or tandem repeat within the sequenced read. Reads with lengths less than twice the given window size are categorized and stored as short reads. Additionally, it incorporates a size selection mechanism to isolate longer reads, thereby streamlining the genome assembly process and enhancing contiguity.

Following evaluation of the CADECT pipeline on fragmented DNA and a comparative analysis of results pre- and post-amplification assay (Table S1), we confirmed that the final genome

8

251    assembly exhibited significantly reduced fragmentation. The integration of MDA and CADECT

252    proved to be effective, particularly in handling challenging, low quantity DNA samples. This

253    combination facilitated the generation of nearly complete genome assemblies with depths

254    above 70× (Fig. 4; Table S4).

255

256



257

258    **Figure 4. Circos plot illustrating a synteny comparison between the reference *S. aureus* ATCC-**

259    **29213 genome sequence and pre- and post-amplification genome assemblies**. The Circos plots

260    contrast the assemblies resulting from the amplified fragmented DNA before and after CADECT

261    processing. A comparison between the reference genome and (A) genomic assembly of

262    fragmented DNA sample without MDA, (B) genomic assembly of fragmented DNA sample with

263    MDA before CADECT and (C) genomic assembly of fragmented DNA sample with MDA after

264    CADECT.

265

266    Overall, when comparing the data before and after CADECT using default parameters with a 500

267    base window size, we observe that its stringent process, which separates putative concatemers

268    and shorter reads, tends to affect the average final depth of the final input. Specifically, in the

269    case of *S. aureus*, we note that for high-quality intact amplified DNA, the detection of putative

270    chimeras and size selection decreases coverage by 40%, whereas for amplified fragmented

271    samples, it decreases coverage by 50% (Table S5). The effect on depth is more pronounced for

272    fragmented samples due to size selection.


273    **DISCUSSION**

274

275    Our study demonstrates that MDA offers a promising solution for amplifying low amounts of DNA

276    of precious samples for ONT sequence generation with the ONT rapid barcode sequencing Library

277    kit (RBK). In this study, we demonstrated three key points: (A) Using our method, we can

278    successfully sequence samples with DNA inputs as low as 0.025 pg. This suggests that long-read

279  sequencing of single cells may be possible. Single-cell sequencing represents a significant
280  achievement as variability, if it exists in the sample, is reduced in the sequence results because
281  we are targeting a considerably smaller number of cells compared to larger bulk-extracted
282  samples. (B) For *Cryptosporidium*, we show that we can reach single-oocyst levels, as 0.025 pg is
283  equivalent to ~1.5 times the amount of DNA in one oocyst (Table 2). This is significant for
284  *Cryptosporidium* because this parasite cannot be cloned.  The prospect of Single-oocyst
285  sequencing removes the variation introduced with bulk sequencing approaches.  (C) We explored
286  and showed the potential of Whole Genome Amplification (WGA) as an option to examine even
287  smaller quantities of DNA depending on the organism under investigation and the size of its
288  genome. Here, we have only examined organisms with genome sizes < 10Mb.  Larger genome
289  sizes will require additional starting material and smaller genome sizes should have success with
290  even less input DNA.

291

292  **Table 2. Estimated DNA concentration in a single cell of the organisms studied in this**
293  **project.**

294

| SAMPLE | Estimated genome size (Mb) | Amt of DNA/cell (fg) |
|---|---|---|
| *E. coli* | 5.00 | 5.43 |
| *S. aureus* | 2.81 | 3.03 |
| *E. faecium* | 2.91 | 3.14 |
| *Cryptosporidium*. ssp. (oocysts)* | 9.2 | 39.70 |
| *Cryptosporidium* ssp. (sporozoite) | | 9.90 |

306  *One *Cryptosporidium* ssp. oocyst contains 4 haploid sporozoites.

307

308  We were able to obtain whole genome sequences at the chromosomal level for almost all tested
309  organisms when generating depth coverages >70×. This indicates that we can overcome the
310  challenge of the relatively shorter reads that MDA with T7 debranching generates, in comparison
311  to reads generated from higher DNA input without amplification. It's essential to highlight that
312  highly complex regions, such as repetitive regions with tandem repeats larger than the window
313  size used for CADECT detection, might be identified as concatemeric reads. This occurs because
314  the tool detects repeat overlaps, which can lead to their exclusion before the assembly process,
315  potentially causing some regions of the genome to remain fragmented. This outcome may vary
316  depending on the organism being sequenced.

317

318    At a lower amount of initial DNA, we observed that the amplification appears to be random,
319    exhibiting no apparent bias across the genome (Fig. 1). At extremely low DNA amounts, achieving
320    full coverage of the target genome may be challenging, but the method remains valuable for
321    potential taxon identification and may prove effective for the identification of plasmids as well
322    (Fig. 1). While the effectiveness in metagenomic samples requires further evaluation, there is
323    promise in using this approach for taxon identification. Extremely low-abundance samples tend
324    to produce patchy sequence information.  Thus, although extremely low-concentration samples
325    provide valuable sequence information, they also lack coverage in many regions, which impacts
326    the assembly process and the ability to produce full genomic information.  Combining multiple
327    MDA replicates is likely to increase the chances of amplifying more regions and thus will be more
328    likely to enhance genome coverage verses deeper sequencing. Interestingly, GC ratios apparently
329    did not impact the amplification showing very low correlations (Fig. S1).

330

331    In the case of sheared DNA, the higher impact on the depth after CADECT is primarily related to
332    loss in the size selection pipeline (Table S5). However, our concatemer detection tool, CADECT,
333    effectively identified and removed several concatemers, facilitating the assembly and yielding
334    good results. This highlights the importance of bioinformatic tools in overcoming challenges
335    associated with amplification artifacts thus improving the accuracy of genome assembly.

336

337    It is worth noting that the CADECT pipeline will remove a significant number of reads which will,
338    impact depth for an optimal genome assembly. If there isn't sufficient coverage obtained post-
339    CADECT run, an alternative is to merge the reads identified as short by the program with the non-
340    concatemeric reads. As observed previously, the chimeric rate produced by MDA is positively
341    associated with the mean read length (Lu et al. 2023), indicating a decreased likelihood of
342    chimeric reads in this short dataset. Consequently, this dataset is less likely to negatively impact
343    the assembly process. In more complex genome sequences that are rich in repeats, further
344    investigation is required to address these regions effectively and be able to distinguish
345    concatemers from genuine repetitive patterns within the genome. As a solution, the CADECT
346    pipeline generates a separate concatemer fastq file. This file includes putative concatemeric
347    regions as well as true repeats.  For example, highly repetitive genomes such as trypanosomatids
348    with an ~50% genomic repeat content (El-Sayed et al. 2005), CADECT would detect a good
349    number of reads containing real tandem repeats in the genome as putative concatemers, which
350    would result in a higher impact on coverage depth loss and also impact the genome content of
351    the organisms used for assembly. To mitigate this, we recommend incorporating a repeat
352    identification step into the pipeline, such as using RepeatModeler (Flynn et al. 2020) trained with
353    the organism of interest on the putative concatemer generated sequence file from CADECT. This
354    additional step would enhance the recovery of information and data for the subsequent assembly
355    process.

356

357 Moving forward, it is crucial to continue exploring the potential of MDA in various biological

358 contexts and optimize the amplification protocol to minimize biases and errors. Additionally,

359 considering the clinical applications of MDA, further research and development of rapid and

360 reliable sequencing approaches are necessary to unlock its full potential in diagnosing and

361 monitoring infectious diseases and other clinical applications.

362


363 **MATERIALS AND METHODS**

364

365 **Sample Collection and Preparation**

366

367 A 100 ng DNA sample of *Cryptosporidium meleagridis* isolate TU1867 was obtained from BEI

368 Resources (NR-2521) (Manassas, VA). DNA samples from cultured *Staphylococcus aureus* ATCC-

369 29213, *Enterococcus faecium* TX-1330, and *Escherichia coli* strain K12, which were available in

370 our laboratory, were used for testing. The bacterial DNA samples were prepared for downstream

371 processing using the QIAGEN QIAamp DNA Mini Kit with lysozyme for Gram-positive samples and

372 buffer ATL (tissue lysis buffer) for Gram-negatives. To assess sequence integrity, an *S. aureus* DNA

373 sample aliquot was sheared using NEBNext dsDNA Fragmentase for 15 minutes to generate

374 fragments approximately 1000 bp in size. All DNA samples were quality controlled using a

375 TapeStation (Agilent Technologies, Santa Clara, CA) and Qubit (Thermofisher Scientific, Walthma,

376 MA). In addition, we conducted serial dilutions on all samples to assess the limit of detection for

377 amplification in the assay. The dilutions ranged from 2.5E-5 ng to 2.5 ng, allowing us to determine

378 the minimum concentration at which successful amplification could be achieved.

379

380 **Multiple Displacement Amplification (MDA)**

381

382 Prior to whole genome amplification (WGA), the concentration of DNA was obtained using a

383 Qubit fluorometer dsDNA high-sensitivity assay kit (ThermoFisher, Waltham, MA). For the *C.*

384 *meleagridis* DNA, three different amounts were used as input for whole genome amplification

385 (*i.e.*, 2 ng, 5 ng, and 10 ng) in a final volume of 5 μL. For the bacterial samples, MDA was

386 performed on 2.5 ng of fragmented of intact *S. aureus DNA* as well as serial dilutions of DNA from

387 *E. faecium* ranging from 2.5 ng to 2.5E-5 ng. 400 ng of non-amplified DNA was used as an input

388 control for the ONT rapid kit library preparation (Oxford Nanopore Technologies, Oxford, United

389 Kingdom).

390 Whole genome amplification (WGA) was performed using the Qiagen Repli-G kit (CAT #150023,

391 Qiagen, Hilden, Germany), following the manufacturer's instructions. Following this,

392     concentrations of DNA were obtained using a Qubit fluorometer dsDNA high-sensitivity assay kit
393     (ThermoFisher, Waltham, MA).
394     For T7 Endonuclease I debranching, up to 42 µL (*i.e.*, all product from the WGA reaction) of WGA
395     DNA was used as input (Catalog #M0302, New England BioLabs, Ipswich, MA). After scaling up
396     the reaction to accommodate a 42 µL input, all reaction components were added following the
397     manufacturer's guidelines and incubated for 15 minutes at 37°C in a BioRad T100 thermocycler
398     (BioRad, Hercules, CA). The incubated reaction was brought up to a final volume of 50 µL using
399     TE buffer pH 8. AMPure XP beads (CAT# A63880) were prepared ahead of time following the
400     manufacturer's instructions, and 35 µL of beads were added to the reaction and mixed
401     thoroughly. The bead-reaction mixture was placed on a rotator mixer (*e.g.*, Hula mixer) for 10
402     minutes at room temperature. Following this, the bead-reaction mixture was spun down and
403     placed on a magnet until the eluate was clear and colorless. With the bead-reaction mixture on
404     the magnet, the clear supernatant was pipetted off and 200 µL of freshly prepared 70% ethanol
405     was carefully added not to disturb the pellet (i.e., wash step). The wash step was repeated one
406     time, for a total of two washes. After removing the supernatant from the second wash, 49 µL of
407     water was used to resuspend the pellet which was immediately incubated for one minute at 50°C
408     in a BioRad T100 thermocycler followed by five minutes at room temperature. The bead-reaction
409     mixture was placed back on the magnet, and 49 µL of the elute was transferred to a sterile 1.5
410     ml tube. Concentrations of DNA were obtained using a Qubit fluorometer dsDNA high-sensitivity
411     assay kit (ThermoFisher, Waltham, MA).
412
413     **Whole Genome Sequencing and Assembly:**
414
415     Sequencing of the amplified DNA samples was performed using the ONT SQK-RBK110.96 kit for
416     library preparation R9.4 MinION flow cells (Oxford Nanopore Technologies, Oxford, UK). The
417     amplified DNA samples were prepared according to the kit instructions and loaded onto the flow
418     cell for sequencing following manufacturer's instructions. Sequencing was carried out in Mk1B
419     and GridION MK1 devices for 72 hours and the resulted fast5 files were basecalled using guppy
420     v6.3.7 using the high accuracy model (dna_r9.4.1_450bps_hac).
421     Flye 2.9 (Kolmogorov et al. 2019) was used for assembly. For samples with >100X coverage, the
422     "--asm-coverage 100" parameter was used to improve assembly and facilitate the assembler
423     performance. We then used Nextpolish 1.4.1 (Hu et al. 2020) to increase the overall basecall
424     quality of the genome and facilitate further quality control analysis such as Benchmarking
425     Universal Single-Copy Orthologs (BUSCO) scores, and better gene annotation. Illumina
426     sequencing was not used here because the objective of this research was to determine how MDA
427     would affect long-read generation.
428
429     **Putative Concatemer Detection in Intact vs. Fragmented Amplified Samples**

430

431    To examine the impact of fragmentation on MDA products, we treated DNA aliquots with
432    NEBNext dsDNA Fragmentase (CAT# M0348S) at 37°C for 16 minutes, creating fragments
433    between 500-1000 bp. This enzyme-based method induces DNA shearing, generating fragments
434    of specified sizes in a time-dependent manner. The process provides random fragmentation
435    similar to mechanical methods. Both fragmented and non-fragmented (high molecular weight
436    DNA) samples were sequenced as described above.

437

438    The CADECT tool (https://github.com/rpbap/CADECT), was developed in-house and was used for
439    the detection and removal of putative concatemeric chimeric sequences in the ONT amplified
440    reads. CADECT splits all reads into separate files and performs sliding windows with a user-
441    defined preferred size and gap between windows. For ONT amplified reads, a window size of >=
442    500 bp with no overlaps was used (*e.g.*, -w 500 and -s 500). Reads generating less than one
443    window (< 1 kb in size in the 500 bp window example) were skipped, and their IDs were stored
444    in the short.txt output file. Fragment windows from reads with more than two windows were
445    aligned using nucmer from mummer 4 (Marcais et al. 2018), and reads with overlaps were
446    reported in the stats file, with their IDs stored in the concat_ID output file. Statistics including the
447    total number of reads, number of putative concatemers, number of reads with no concatemer
448    detection, and overlap frequency were recorded in the stats.txt output file. Fastq/Fasta files
449    containing the characterized reads were generated for further analysis.

450

451    These methods were employed to investigate the benefits and limitations of multiple
452    displacement amplification in whole-genome Oxford Nanopore Sequencing, focusing on low-
453    concentration DNA samples.

454

455    **GC Bias Evaluation**

456

457    To calculate GC bias in the sequencing depth of the amplified data, we compared the local %GC
458    content using sliding windows of 1000 bp to the average coverage depth for each of these
459    regions (https://github.com/DamienFr/GC_content_in_sliding_window). Depth windows were
460    calculated using the R packages setDT and rollapply packages. $R^2$ coefficients were calculated
461    using the ordinary least squares regression method.

462    **DATA AVAILABILITY**

463

464    The raw sequence data generated in this study have been submitted to the NCBI BioProject
465    database (https://www.ncbi.nlm.nih.gov/bioproject/) under accession numbers PRJNA1063853
466    and PRJNA1022047. The assembled genomes in this project are preliminary drafts and are

14

467  currently unavailable at this stage due to the scope of this project. They were generated solely
468  using ONT reads and have not undergone polishing with Illumina short-read data. Additionally,
469  they have not been checked for potentially contaminating "leftover" contigs. The raw data is
470  accessible for reproduction purposes, and the final, polished, and decontaminated assemblies
471  will be made available in subsequent publications.

472  **COMPETING INTEREST STATEMENT**

473  The authors declare that there are no conflicts of interest.

474  **ACKNOWLEDGMENTS**

475

479

480  **AUTHOR CONTRIBUTIONS**

481

482  FAD, MSB, AD, HS, MQD, RPB performed the research. MSB, JCK, TCG, RPB conceived the
483  research. FAD, MSB, AD, MQD and RPB analyzed the results. FAD, MSB, AD, RT, AK, CAA, JCK, TCG,
484  RPB contributed to writing the manuscript. CAA, JCK, TCG and RPB obtained funding. All authors
485  reviewed and approved of the manuscript.
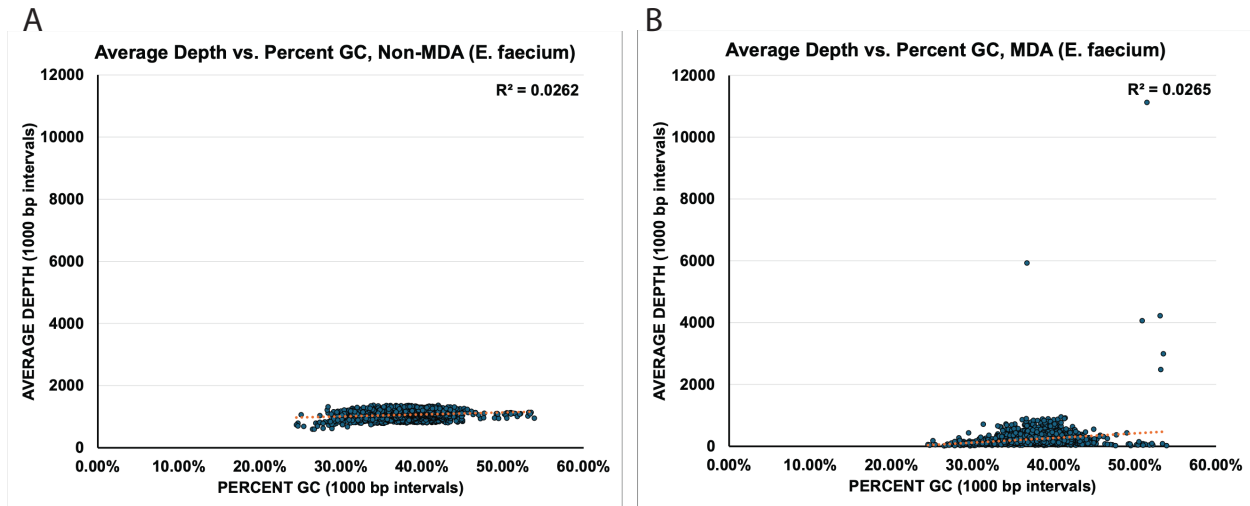
486

487

488

## REFERENCES

490

491 Binga EK, Lasken RS, Neufeld JD. 2008. Something from (almost) nothing: the impact of
492       multiple displacement amplification on microbial ecology. *ISME J* **2**: 233-241.
493 Ceccherini M, Pote J, Kay E, Van VT, Marechal J, Pietramellara G, Nannipieri P, Vogel TM,
494       Simonet P. 2003. Degradation and transformability of DNA from transgenic leaves. *Appl*
495       *Environ Microbiol* **69**: 673-678.
496 Dean FB, Hosono S, Fang L, Wu X, Faruqi AF, Bray-Ward P, Sun Z, Zong Q, Du Y, Du J et al.
497       2002. Comprehensive human genome amplification using multiple displacement
498       amplification. *Proc Natl Acad Sci U S A* **99**: 5261-5266.
499 Delahaye C, Nicolas J. 2021. Sequencing DNA with nanopores: Troubles and biases. *PLoS*
500       *One* **16**: e0257521.
501 El-Sayed NM, Myler PJ, Bartholomeu DC, Nilsson D, Aggarwal G, Tran AN, Ghedin E, Worthey
502       EA, Delcher AL, Blandin G et al. 2005. The genome sequence of *Trypanosoma cruzi*,
503       etiologic agent of Chagas disease. *Science* **309**: 409-415.
504 Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. 2020.
505       RepeatModeler2 for automated genomic discovery of transposable element families.
506       *Proc Natl Acad Sci U S A* **117**: 9451-9457.
507 Fullwood MJ, Tan JJ, Ng PW, Chiu KP, Liu J, Wei CL, Ruan Y. 2008. The use of multiple
508       displacement amplification to amplify complex DNA libraries. *Nucleic Acids Res* **36**: e32.
509 Hard J, Mold JE, Eisfeldt J, Tellgren-Roth C, Haggqvist S, Bunikis I, Contreras-Lopez O, Chin
510       CS, Nordlund J, Rubin CJ et al. 2023. Long-read whole-genome analysis of human
511       single cells. *Nat Commun* **14**: 5164.
512 Hou Y, Wu K, Shi X, Li F, Song L, Wu H, Dean M, Li G, Tsang S, Jiang R et al. 2015.
513       Comparison of variations detection between whole-genome amplification methods used
514       in single-cell resequencing. *Gigascience* **4**: 37.
515 Hu J, Fan J, Sun Z, Liu S. 2020. NextPolish: a fast and efficient genome polishing tool for long-
516       read assembly. *Bioinformatics* **36**: 2253-2255.
517 Hu T, Chitnis N, Monos D, Dinh A. 2021. Next-generation sequencing technologies: An
518       overview. *Hum Immunol* **82**: 801-811.
519 Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using
520       repeat graphs. *Nat Biotechnol* **37**: 540-546.
521 Lasken RS. 2009. Genomic DNA amplification by the multiple displacement amplification (MDA)
522       method. *Biochem Soc Trans* **37**: 450-453.
523 Lasken RS, Stockwell TB. 2007. Mechanism of chimera formation during the Multiple
524       Displacement Amplification reaction. *BMC Biotechnol* **7**: 19.
525 Lu N, Qiao Y, Lu Z, Tu J. 2023. Chimera: The spoiler in multiple displacement amplification.
526       *Comput Struct Biotechnol J* **21**: 1688-1696.
527 Marcais G, Delcher AL, Phillippy AM, Coston R, Salzberg SL, Zimin A. 2018. MUMmer4: A fast
528       and versatile genome alignment system. *PLoS Comput Biol* **14**: e1005944.

529    Pacbio. 2022. Technical note Preparing DNA for PacBio HiFi sequencing—extraction and
530         quality control.
531    Paul P, Apgar J. 2005. Single-molecule dilution and multiple displacement amplification for
532         molecular haplotyping. *Biotechniques* **38**: 553-554, 556, 558-559.
533    Slatko BE, Gardner AF, Ausubel FM. 2018. Overview of Next-Generation Sequencing
534         Technologies. *Curr Protoc Mol Biol* **122**: e59.
535    Wang G, Maher E, Brennan C, Chin L, Leo C, Kaur M, Zhu P, Rook M, Wolfe JL, Makrigiorgos
536         GM. 2004. DNA amplification method tolerant to sample degradation. *Genome Res* **14**:
537         2357-2366.
538
539
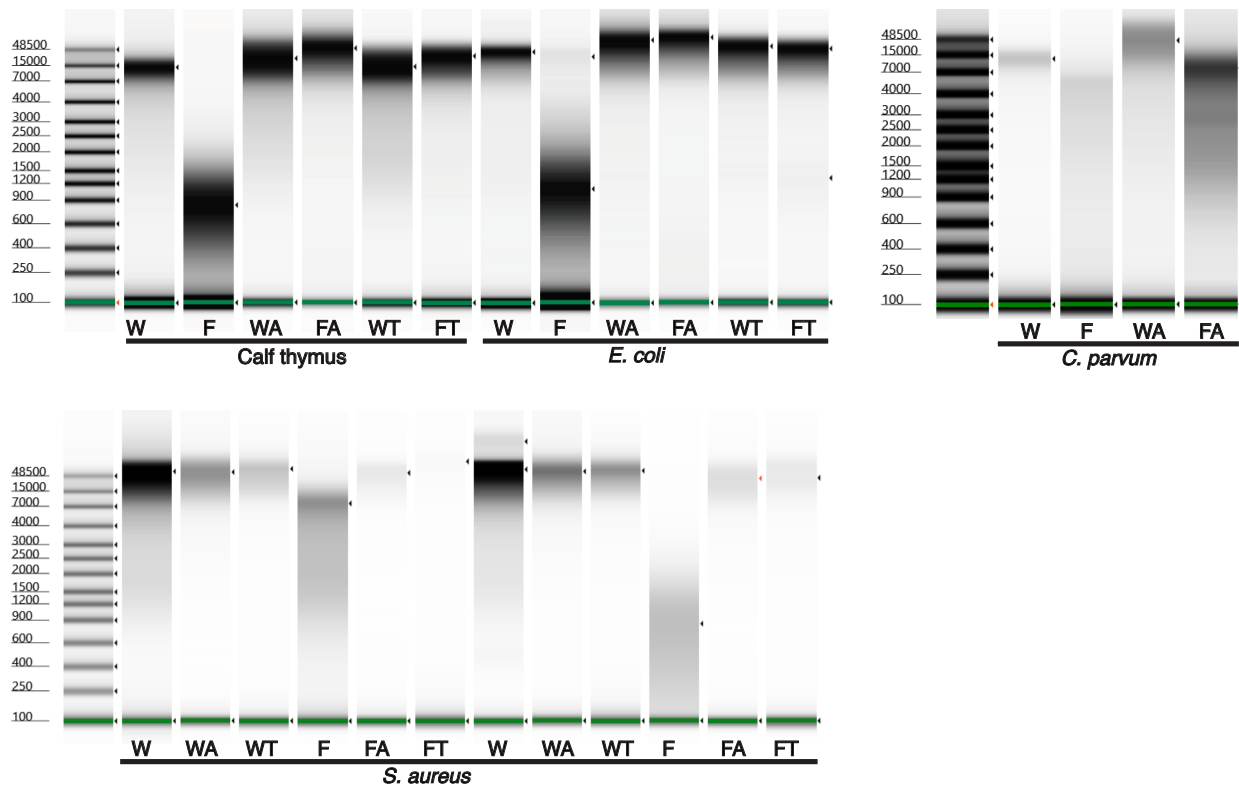
540 **SUPPLEMENTAL FIGURES**

541

542



543

544

545 **Figure S1**. **Square regression correlation coefficient between average depth and average %GC**

546 **across 1000 base pair sliding window regions of the genomic assembly shows low correlation**

547 **to %GC.** A correlation analysis of (A) non-amplified and (B) amplified assemblies of *E. faecium.*

548

549



550
551

**Figure S2.** Uncropped DNA Fragment and read Size Range pre- and post- Multiple Displacement Amplification using Size-Controlled Fragmented DNA. Uncropped TapeStation results for different organism sets; (W = whole intact DNA; F = Fragmented DNA; WA and FA = after amplification; and WT and FT = after T7 debranching.

556
557

**SUPPLEMENTAL TABLES**

**Table S1** - Whole genome assembly statistics for the data generated using MDA.

| Assembly | Condition | Total length (nt) | # of contigs | Largest contig (nt) | GC% | N50 | L50 | # N's per 100 Kbp |
|---|---|---|---|---|---|---|---|---|
| *C. meleagridis* | WGA (5 ng) | 9,171,013 | 13 | 1,363,785 | 30.92 | 1,103,979 | 4 | 0 |
| *E. coli* | WGA-Intact DNA | 6,702,699 | 252 | 479,759 | 50.49 | 157,769 | 12 | 0 |
| | WGA-Fragmented DNA | 157,593 | 12 | 25,757 | 49.23 | 20,338 | 4 | 0 |
| *E.faecium* | 400ng (no MDA) | 2,775,595 | 2 | 2,583,377 | 38.27 | 2,583,377 | 1 | 0 |
| | 1 ng | 3,603,364 | 216 | 257,348 | 38.3 | 64,352 | 14 | 0 |
| | 1E-1ng | 4,287,796 | 570 | 83,451 | 38.66 | 27,070 | 49 | 0 |
| | 1E-2 ng | 4,146,213 | 1,066 | 65,055 | 38.94 | 13,948 | 95 | 0 |
| | 1E-3 ng | 1,174 | 2 | 670 | 53.15 | 670 | 1 | 0 |
| | 1E-4 ng | 78,463 | 53 | 15,260 | 53.83 | 2,204 | 7 | 0 |
| | 1E-5 ng | 54,253 | 32 | 6,962 | 56.49 | 3,816 | 6 | 0 |
| *S. aureus* | Intact DNA | 2,766,204 | 3 | 2,717,982 | 32.86 | 2,717,982 | 1 | 0 |
| | Fragmented DNA | 854,501 | 198 | 33,734 | 33.62 | 4,926 | 36 | 0 |
| | Intact DNA post-WGA | 2,763,611 | 4 | 2,717,354 | 32.86 | 2,717,354 | 1 | 0 |
| | Fragmented DNA post-WGA | 2,842,696 | 20 | 1,890,364 | 32.84 | 1,890,364 | 1 | 0 |

20

563 **Table S2** - Sequencing Statistics for *E. faecium\** WGS results at 2.5 ng and 0.2.5 ng starting DNA
564 input.
565

| Starting total DNA | 2.5 ng | 0.25 ng |
|---|---|---|
| Depth | 73× | 35× |
| Total length (bp) | 2,778,112 | 2,918,507 |
| Number of contigs | 3 | 35 |
| GC% | 38.2 | 38.28 |
| N50 | 2,589,111 | 233,792 |
| L50 | 1 | 4 |
| Number of N's per 100 Kbp | 0 | 0 |

566 *the expected genome size for *E. faecium* species varies between 2.6 and 3.2 Mb
567

568    **Table S3** - Read length distribution among the DNA dilutions for *E. faecium,* based on starting

569    DNA concentration.

570

| Starting DNA Amount (ng) | Mean Read Length (bp) | Median Read Length (bp) |
|---|---|---|
| Control 400 | 6,134 | 3,819 |
| 2.5 | 2,533 | 1,308 |
| 2.5E-01 | 2,305 | 1,002 |
| 2.5E-02 | 2,390 | 1,147 |
| 2.5E-03 | 3,121 | 1,086 |
| 2.5E-04 | 2,555 | 1,185 |
| 2.5E-05 | 2,641 | 957 |

571

572

573 **Table S4.** Comparison between *S. aureus* ATCC-29213 genome sequences assembled before

574 and after CADECT with assembly length, number of contigs, number of reads, and mean read

575 length.

576

|  | Fragmented DNA with MDA <u>before</u> CADECT | Fragmented DNA with MDA <u>after</u> CADECT |
|---|---|---|
| **Assembly length (bp)** | 2,842,696 | 2,752,482 |
| **Number of contigs** | 20 | 3 |
| **Number of reads** | 324,197 | 292,789 |
| **Mean read length (bp)** | 2,163 | 1,534 |
| **Median read length (bp)** | 1,352 | 1,202 |

577

578 **Table S5.** Comparison between *S. aureus* results from CADECT using low DNA input samples.

579

| DNA input Condition | Intact DNA | | Fragmented DNA | |
|---|---|---|---|---|
| MDA | No | yes | no | yes |
| **Sequencing coverage input** | 46.8 | 92.3 | 21.7**** | 77.4 |
| **Total number of sequenced base pairs** | 131,012,134 | 258,412,289 | 60,893,145 | 216,785,548 |
| **Number of shorter reads detected*** | 6,397 | 49,746 | 109,830 | 54,268 |
| **Number of non-concatemer reads detected** | 12,556 | 44,440 | 2,548 | 35,986 |
| **Number of putative concatemer reads detected** | 1,047 | 5,814 | 26** | 9,746** |
| **Total number of Reads analyzed:** | 20,000 | 100,000 | 112,404 | 100,000 |
| **read loss (%)** | 37 | 56 | 98 | 64 |
| **Total number of non-Concatemer base pairs** | 103,146,375 | 153,374,280 | 5,304,400 | 96,993,865 |
| **Coverage loss (x)** | 10.0 | 37.5 | 19.9 | 42.8 |

580 *Shorter reads were reads detected below the default setting of 500 nt window size

581 **Due to size selection putative concatemers were classified as short reads

582 ***Loss if using just the reads characterized as non-concatemeric

583 ****Without the amplification ONT had a bad throughput for the fragmented samples at low

584 input values to generate longer reads, resulting in a low sequencing coverage.