

1 **Title:** Long-read assembly and comparative evidence-based reanalysis of *Cryptosporidium*  
2 genome sequences reveal new biological insights  
3

4 **Authors:** Rodrigo P. Baptista<sup>1,2</sup>, Yiran Li<sup>2</sup>, Adam Sateriale<sup>3</sup>, Mandy J. Sanders<sup>4</sup>, Karen L. Brooks<sup>4</sup>,  
5 Alan Tracey<sup>4</sup>, Brendan R. E. Ansell<sup>5</sup>, Aaron R. Jex<sup>5</sup>, Garrett W. Cooper<sup>6</sup>, Ethan D. Smith<sup>6</sup>, Rui  
6 Xiao<sup>2</sup>, Jennifer E. Dumaine<sup>3</sup>, Matthew Berriman<sup>4</sup>, Boris Striepen<sup>3</sup>, James A. Cotton<sup>4</sup> and Jessica  
7 C. Kissinger<sup>1,2,6</sup>  
8

9 **Affiliation:** <sup>1</sup>Center for Tropical and Emerging Global Diseases; <sup>2</sup>Institute of Bioinformatics,  
10 University of Georgia, Athens, GA, USA; <sup>3</sup>Department of Pathology, School of Veterinary Medicine,  
11 University of Pennsylvania, Philadelphia, PA, USA; <sup>4</sup>The Wellcome Sanger Institute, Hinxton, UK;  
12 <sup>5</sup>Faculty of Veterinary and Agricultural Sciences, The University of Melbourne and Population  
13 Health and Immunity Division, the Walter and Eliza Hall Institute of Medical Research, Melbourne,  
14 Australia and <sup>6</sup>Department of Genetics, University of Georgia, Athens, GA, USA  
15

16 **Key words:** Reassembly, Long-read, genome rearrangement, telomere replication,  
17 subtelomere, DNA methylation, reannotation  
18

19 **Corresponding Author:** Jessica C. Kissinger [jkissing@uga.edu](mailto:jkissing@uga.edu)  
20

21 **Manuscript Type:** Research  
22

23 **Running Title:** New Comparative *Cryptosporidium* Genomic Insights  
24

## 25 **ABSTRACT**

26  
27 Cryptosporidiosis is a leading cause of waterborne diarrheal disease globally and an important  
28 contributor to mortality in infants and the immunosuppressed. Despite its importance, the  
29 *Cryptosporidium* community still relies on a fragmented reference genome sequence from 2004.  
30 Incomplete reference sequences hamper experimental design and interpretation. We have  
31 generated a new *C. parvum* IOWA genome assembly supported by PacBio and Oxford Nanopore  
32 long-read technologies and a new comparative and consistent genome annotation for three  
33 closely related species *C. parvum*, *C. hominis* and *C. tyzzeri*. The new *C. parvum* IOWA reference  
34 genome assembly is larger, gap free and lacks ambiguous bases. This chromosomal assembly  
35 recovers 13 of 16 possible telomeres and raises a new hypothesis for the remaining telomeres  
36 and associated subtelomeric regions. Comparative annotation revealed that most “missing”  
37 orthologs are found suggesting that species differences result primarily from structural  
38 rearrangements, gene copy number variation and SNVs in *C. parvum*, *C. hominis* and *C. tyzzeri*.  
39 We made >1,500 *C. parvum* annotation updates based on experimental evidence. They included  
40 new transporters, ncRNAs, introns and altered gene structures. The new assembly and  
41 annotation revealed a complete DNA methylase *Dnmt2* ortholog. 190 genes under positive  
42 selection including many new candidates were identified using the new assembly and annotation  
43 as reference. Finally, possible subtelomeric amplification and variation events in *C. parvum* are  
44 detected that reveal a new level of genome plasticity that will both inform and impact future  
45 research.  
46  
47  
48  
49

## 50 **INTRODUCTION**

51

52 *Cryptosporidium* spp. are parasitic apicomplexans that cause moderate-to-severe  
53 diarrhea in humans and animals. Studies funded by the Bill and Melinda Gates Foundation,  
54 revealed that *Cryptosporidium* is one of the most common causes of waterborne disease in  
55 humans and the second leading cause of diarrheal etiology in children < 2 years resulting in  
56 ~60,000 fatalities worldwide (Kotloff et al. 2013; Collaborators 2017). In 2016, acute infections  
57 caused more than 48,000 global deaths and more than 4.2 million disability-adjusted life years  
58 lost (Khalil et al. 2018).

59 Currently, 38 species of *Cryptosporidium* are recognized by the scientific community  
60 (Slapeta 2013; Feng et al. 2018). Most are host-adapted, and host species range from fish to  
61 mammals. Of these, 15 species have had their genome sequence generated and assembled  
62 however, only 8 are annotated. Most genomic sequence data are from the zoonotic *C. parvum*  
63 and anthroponotic *C. hominis*, the species primarily detected in humans (Chalmers et al. 2011;  
64 Zahedi et al. 2016; Khan et al. 2017). These two species are only 3-5% divergent at the DNA level  
65 (Mazurie et al. 2013).

66 As the *Cryptosporidium* field is exploding with new-found interest and much needed  
67 breakthroughs in genetics and culturing (Vinayak et al. 2015; Morada et al. 2016; DeCicco  
68 RePass et al. 2017; Heo et al. 2018; Wilke et al. 2019), the limitations of existing reference  
69 genome sequences need to be addressed. The *C. parvum* IOWA II reference genome sequence  
70 was assembled with limited physical map data (Abrahamsen et al. 2004) and experimental data  
71 for training gene finders and providing functional annotation were limited to a few hundred ESTs  
72 from oocysts and sporozoite stages only. Genomic, transcriptomic and proteomic work on this  
73 important pathogen has been lacking due to the obligate quasi-intracellular nature of portions of  
74 the parasite's life cycle, the historical lack of a continuous *in vitro* tissue culture system, the  
75 parasite's small size relative to host cells and difficult animal models. The physical map for the *C.*  
76 *parvum* IOWA II reference assembly was generated from two different studies that utilized the  
77 genome-wide HAPpilly anchored physical mapping technique, an *in vitro* linkage technique based  
78 on screening approximately haploid amounts of DNA by PCR, which is very accurate (Piper et al.  
79 1998; Bankier et al. 2003). Even with these cutting-edge approaches at the time, some regions,  
80 especially chromosome ends, lacked support or were poorly resolved. Subsequent whole genome  
81 sequencing data often remain unassembled or in a large number of contigs.

82 In 2015, the reference genome sequence of *C. parvum* was re-annotated based on new  
83 RNA-seq evidence and a new *C. hominis* sequence from a recent human isolate (UdeA01) was  
84 generated (Isaza et al. 2015). Many ambiguities in gene models were improved based on the new  
85 RNA-seq data, but since the new *C. hominis* UdeA01 genome is still fragmented and the  
86 annotation was primarily based on the 2004 *C. parvum* IOWA II reference annotation. Additionally,  
87 annotation of sequences from closely related species has been performed independently and are  
88 not consistent, causing possible misinterpretations regarding gene content and species-specific  
89 genes.

90 Incomplete and misassembled (i.e. gapped sequence, indels, frameshifts, compressed  
91 repetitive regions, inversions) reference genome sequences such as those shown in (Guo et al.  
92 2015) can mislead interpretations of the differences between isolates and species resulting in  
93 extra assays to confirm insertions, deletions and copy number variations (CNVs). Since  
94 incomplete and misassembled sequences are usually caused by repetitive and complex  
95 sequence regions, it is imperative to revisit older reference genome sequences with new long-  
96 read technologies to close gaps and expand regions of the genome sequence that were  
97 misassembled or collapsed into shorter regions because they are repetitive. Long-read sequence  
98 technologies (PacBio and Oxford Nanopore) are becoming an essential tool to close full genome  
99 sequence assemblies across the tree of life (Vembar et al. 2016; Diaz-Viraque et al. 2019; Miga  
100 et al. 2020). They can be used to resolve complex regions such as repetitive content, structural  
101 variants (SVs) such as inversions, translocations and duplications, or for use as scaffolding  
102 evidence for existing fragmented genome assemblies (Mahmoud et al. 2019). They are proving

103 crucial for completing assemblies of pathogen genome sequences that are often riddled with large  
 104 virulence-related gene families that have been collapsed or improperly assembled in shorter-read  
 105 assemblies. Here we provide a new *de novo* reference long-read assembly for *C. parvum* strain  
 106 IOWA (DNA obtained from the ATCC) and new consistent, comparative genome annotations for  
 107 *C. parvum* IOWA-ATCC, *C. hominis* UdeA01 and *C. tyzzeri* UGA55.

## 108 RESULTS

### 109 An improved long-read based genome assembly for *Cryptosporidium parvum* (IOWA-ATCC)

110  
 111 The current *C. parvum* IOWA II reference genome assembly, generated in 2004, is good,  
 112 but it still has 10 gapped regions of unknown size, 14,600 ambiguous bases, and is missing 6  
 113 telomeres. By aligning Illumina reads against this reference sequence, we have detected many  
 114 collapsed regions, suggesting misassembled repetitive and complex regions (Supplemental Table  
 115 S1). To resolve these issues, we generated a new PacBio+Illumina+Nanopore hybrid genome  
 116 assembly for the *Cryptosporidium parvum* strain IOWA (ATCC®PRA-67DQ™) with DNA from  
 117 oocysts/sporozoites purchased from ATCC. To minimize strain variation differences, we  
 118 performed our analysis on the same strain, however because there is a 14-year time window of  
 119 propagation between these two isolates, and cryopreservation has only been recently made  
 120 possible (Jaskiewicz et al. 2018), we modified the strain name to IOWA-ATCC.

121 The new *C. parvum* IOWA-ATCC genome statistics are compared to the current *C. parvum*  
 122 IOWA II reference genome sequence and *C. hominis* 30976 and *C. tyzzeri* UGA55 two closely  
 123 related species with different host preferences and pathogenicity (Slapeta 2013; Nader et al. 2019;  
 124 Sateriale et al. 2019) (Table 1). These particular *C. hominis* and *C. tyzzeri* assemblies were  
 125 selected because they are the best available. The new long-read assembly increases the genome  
 126 size by 19,939 bases and identifies 13 of 16 expected telomeres. There are no gaps and no  
 127 ambiguous bases. As expected, the *C. parvum* IOWA-ATCC genome sequence has diverged  
 128 slightly but shares 99.93% average pairwise identity with the 2004 assembly in regions that exist  
 129 in both assemblies (Supplemental Table S2). The main *Cryptosporidium* subtyping marker, the 60  
 130 kDa surface protein (*gp60* locus subtype IIa) shows 4 amino acid differences between the IOWA-  
 131 ATCC and 2004 assemblies (Supplemental Fig. S1).

132  
 133  
 134  
 135

**Table 1 *Cryptosporidium* Genome Assembly Statistics**

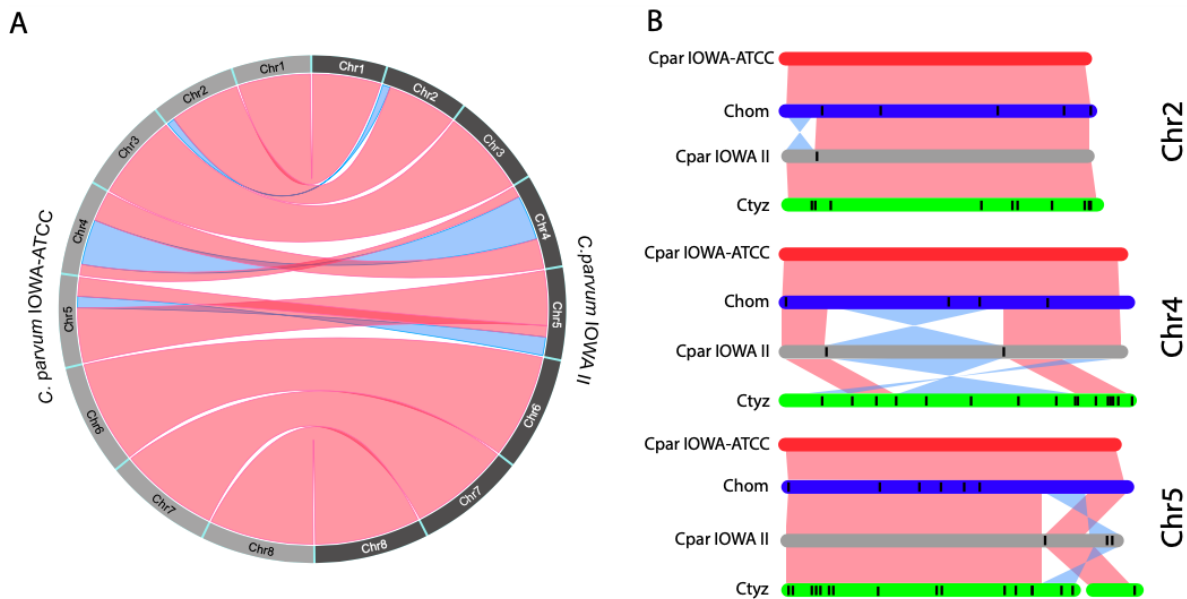
|                            | <i>C. parvum</i><br>IOWA II (2004) | <i>C. parvum</i> IOWA-<br>ATCC | <i>C. hominis</i><br>30976 | <i>C. tyzzeri</i><br>UGA55 |
|----------------------------|------------------------------------|--------------------------------|----------------------------|----------------------------|
| <b>Scaffolds</b>           | 8                                  | 8                              | 53                         | 11                         |
| <b>Gaps in assembly</b>    | 10                                 | 0                              | 25                         | 97                         |
| <b>Total length bp</b>     | 9,102,324                          | 9,122,263                      | 9,059,225                  | 9,015,713                  |
| <b>Compressed regions*</b> | > 14                               | > 8                            | > 18                       | > 17                       |
| <b>Ambiguous nt "N's"</b>  | 14,600                             | 0                              | 1,699                      | 78,408                     |
| <b># Of telomeres</b>      | 10                                 | 13                             | 7                          | 8                          |
| <b>N50</b>                 | 1,104,417                          | 1,108,396                      | 470,636                    | 1,108,290                  |
| <b>GC%</b>                 | 30.23                              | 30.18                          | 30.13                      | 30.25                      |

136 \*Numbers represents compressed regions of > 100nt length and > 3 copies as average depth.

137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156

### Structural differences between the *C. parvum* IOWA assemblies

The 2004 *C. parvum* IOWA II genome assembly used Sanger reads combined with available HAPPY-map data to scaffold the contigs. We compared the 2004 and IOWA-ATCC assemblies to identify potential rearrangements. Small and large rearrangements were detected primarily in chromosomes 2, 4 and 5 (Fig. 1A). Chromosomal inversions may be assembly artifacts or represent genuine differences generated during evolution. Inversions are often associated with speciation events (de Meeus et al. 1998; Rieseberg 2001; Nosil and Feder 2012). We thus investigated the synteny between *C. parvum* IOWA II and ATCC, *C. hominis* 30976 and *C. tyzzeri* UGA55 and observed that *C. hominis* and *C. tyzzeri* also share the large inversions in their chr 4 and chr 5. Examination of the inverted region boundaries revealed that sequences in these regions in the 2004 *C. parvum* assembly consist of ambiguous nucleotide bases or physical gaps (Fig. 1B). These results suggest that the 2004 *C. parvum* assembly may contain misassembled scaffolds, but the data do not rule out their presence in that isolate. Better assemblies will be needed for the other isolates to determine the true level of synteny across these species.



**Figure 1.** Syntenic relationships and inversions detected between the *Cryptosporidium* genome assemblies. (A) Circos plot of synteny between *C. parvum* IOWA-ATCC and IOWA II. (B) Synteny between chromosomes 2, 4 and 5 of *C. hominis* 30976, *C. parvum* IOWA II and *C. tyzzeri* UGA55. Each vertical black line within a chromosome represents a known gap region. Syntenic regions between chromosomes are shown in red and inverted regions are represented in blue. Cpar: *C. parvum*; Chom: *C. hominis*; Ctyz: *C. tyzzeri*.

157  
158  
159  
160  
161  
162

### New consistent annotation across *Cryptosporidium* species provides insights

We consistently annotated and compared the three closely related, yet biologically different, *Cryptosporidium* species (genome identity > 95%) to assess differences in gene content. The new annotation for each species was generated with three *de novo* approaches and evidence-based manual annotation. Curation of the annotation was performed in a 3-way

163 comparison between each pair of genome sequences to take full advantage of syntenic regions.  
 164 The comparison permitted the use of data from one species to assess computational predictions  
 165 in the others. By following this approach, fragments of genes that were previously missed in *C.*  
 166 *hominis* were identified, permitting a more accurate identification of genuinely shared and  
 167 species-specific genes in these species. This approach resulted in > 1500 gene structure  
 168 alterations leading to an improved functional annotation. The changes increased the number of  
 169 predicted genes, introns and exons (Table 2). The average mRNA length increased due to  
 170 complete coding sequences (CDS) and the addition of exons to form larger genes. Notably, these  
 171 structural fixes led to the repair of several genes, including finding and correcting the N-terminus  
 172 of the DNA methylase ortholog, *Dnmt2* (Supplemental Fig. S2).

173 *Cryptosporidium* has a very compact genome with < 20% being intergenic. As a result,  
 174 RNA-Seq data, which is the best evidence for annotation, contains reads that overlap adjacent  
 175 genes creating false fusions of exons belonging to different genes. Available strand-specific RNA-  
 176 seq was used to characterize some of these regions but expression data were not available for  
 177 all predicted genes, thus, genes of unknown function in close proximity on the same strand remain  
 178 problematic. The expression data also revealed alternative splicing and potential non-coding  
 179 RNAs (ncRNAs) predominantly anti-sense lncRNAs with differential expression (Li et al. 2020).  
 180

**Table 2 - Reannotation Summary Statistics.**

|                            | <i>C. parvum</i> IOWA II |                  |           | <i>C. hominis</i>   |                  | <i>C. tyzzeri</i><br>UGA55 |
|----------------------------|--------------------------|------------------|-----------|---------------------|------------------|----------------------------|
|                            | IOWA II<br>Before        | IOWA II<br>After | IOWA-ATCC | UdeA01*<br>"Before" | 30976<br>"After" | New                        |
| Total sequence length (bp) | 9,102,324                | 9,102,324        | 9,122,263 | 9,043,938           | 9,059,225        | 9,015,884                  |
| Number of genes            | 3,886                    | 4,020            | 3,954     | 3,863               | 3,996            | 4,037                      |
| Number of CDS              | 3,805                    | 3,944            | 3,944     | 3,818               | 3,959            | 3,986                      |
| Number of exons            | 4,104                    | 5,043            | 5,322     | 4,546               | 5,045            | 5,136                      |
| Number of introns          | 238                      | 1,020            | 1,371     | 683                 | 1,040            | 1,089                      |
| Shortest intron (bp)       | 9                        | 36               | 36        | 36                  | 36               | 22                         |
| Pseudogenes                | 74                       | 114              | 1         | 45                  | 88               | 62                         |
| % of genome covered by CDS | 75.4                     | 82.1             | 79.53     | 76.1                | 83.6             | 79.2                       |

181 \*A previous annotation does not exist for the *C. hominis* 30976 sequence, so the results are  
 182 compared to the recently annotated *C. hominins* UdeA01 strain.  
 183

## 184 185 Functional annotation

186  
 187 Several approaches to assess function were applied including InterPro scan and I-  
 188 TASSER among others (see methods). 138 new protein annotations were generated or modified,  
 189 the rest are unchanged. The percentage of *C. parvum* genes annotated as uncharacterized  
 190 proteins was reduced from 40% to 33% in all reannotated sequences (Supplemental Table S3).  
 191 Many new features including domain and repeat content were added to 738 previously  
 192 uncharacterized proteins. 729 predicted *C. parvum* CDSs have signal peptides and 1990 have

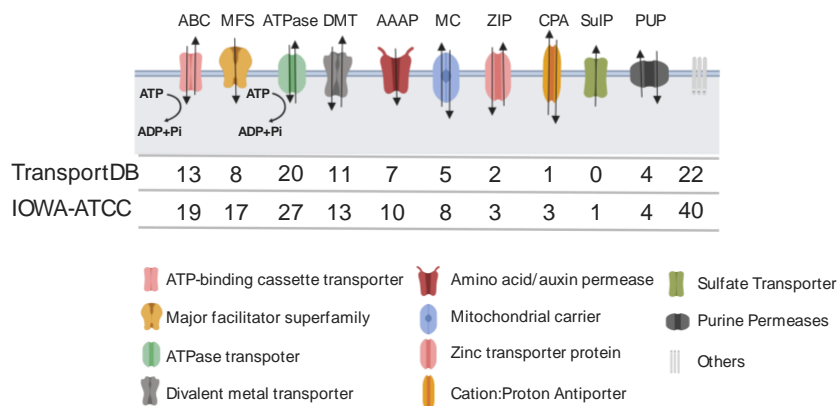
193 GO assignments. 1414 CDSs were further assessed for confidence using I-TASSER protein  
 194 structure searches and 1008 predicted structures were assigned as high-confidence by random  
 195 forest categorization (Supplemental Table S4). 143 previously uncharacterized proteins were  
 196 assigned with high confidence GO terms. The top functional annotation terms observed following  
 197 re-annotation were protein kinases, AAA+ATPases, TRAP, DEAD/DEAH box proteins, Ras  
 198 GTPases, WD40-repeat containing proteins, ABC transporters, RNA recognition motifs,  
 199 Palmitoyltransferases and insulinase-like proteases.

200  
 201

## 202 New transporters were detected and annotated

203

204 Following functional annotation, we further characterized the newly identified transporter  
 205 genes using three different prediction methods. A total of 145 proteins in *C. parvum* IOWA-ATCC  
 206 and *C. hominis* 30976 were identified as transporters including 128 confident candidates and 24  
 207 putative candidates (Supplemental Table S5). This represents an increase of 53 transporters  
 208 relative to the *C. parvum* IOWA II GO annotation (CryptoDB Release 36) and an increase of 93  
 209 relative to TransportDB v2.0 (<http://www.membranetransport.org/transportDB2/index.html>). The  
 210 predicted transporters in *Cryptosporidium* are mostly related to purine metabolism, peptidoglycan  
 211 biosynthesis, oxidative phosphorylation and N-Glycan biosynthesis pathways (Fig. 2). Six  
 212 translocases were also identified.  
 213



214  
 215  
 216

217 **Figure 2.** Reannotation reveals new transporters in *Cryptosporidium parvum* IOWA-ATCC.  
 218 Numbers of transporters corresponds to the counts of genes encoding each type of transporter  
 219 protein. ABC: ATP-binding cassette transporter; MFS: Major facilitator superfamily; DMT: Divalent  
 220 metal transporter; AAAP: amino acid/auxin permease; MC: mitochondrial carrier; ZIP: Zinc  
 221 transporter protein; CPA: Cation/Proton Antiporter; SulP: Sulfate Transporter; and PUP: Purine  
 222 Permeases.

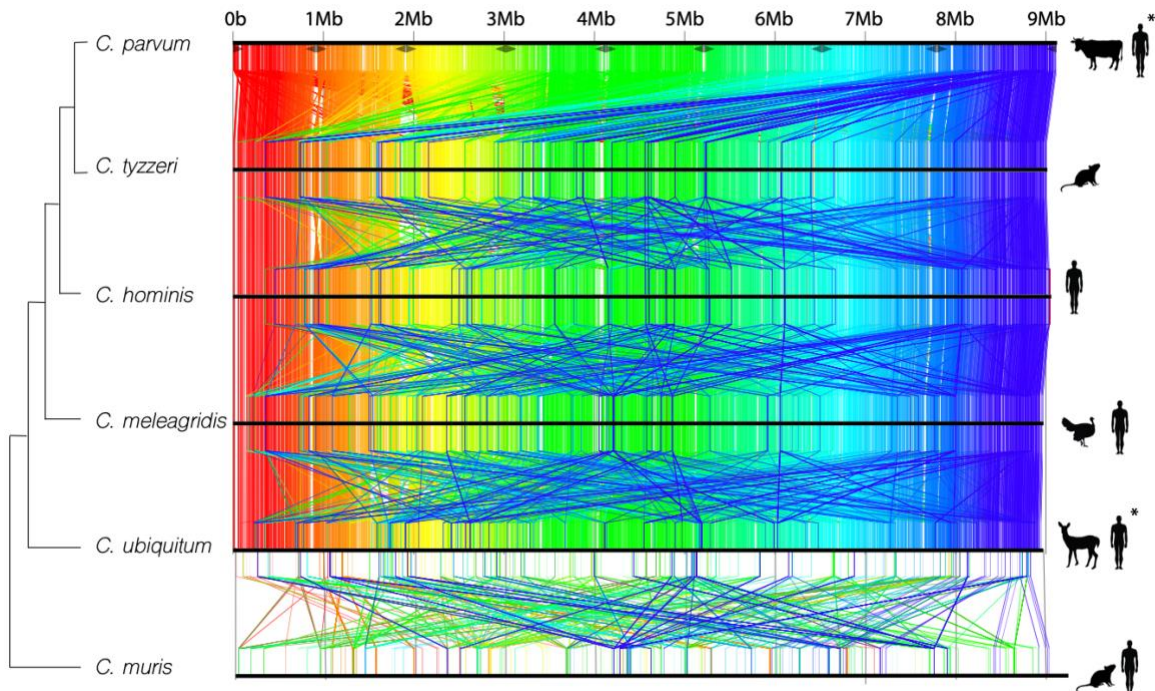
223

## 224 Comparative analysis of closely related species of *Cryptosporidium*

225

226 *Cryptosporidium* species have a broad host spectrum with most species being largely  
 227 host-adapted with a few zoonotic exceptions, principally *C. parvum*. Yet, despite these differences,  
 228 there is a cluster of species with high synteny relative to other species outside of this cluster (Fig.  
 229 3; Supplemental Table S6).

230



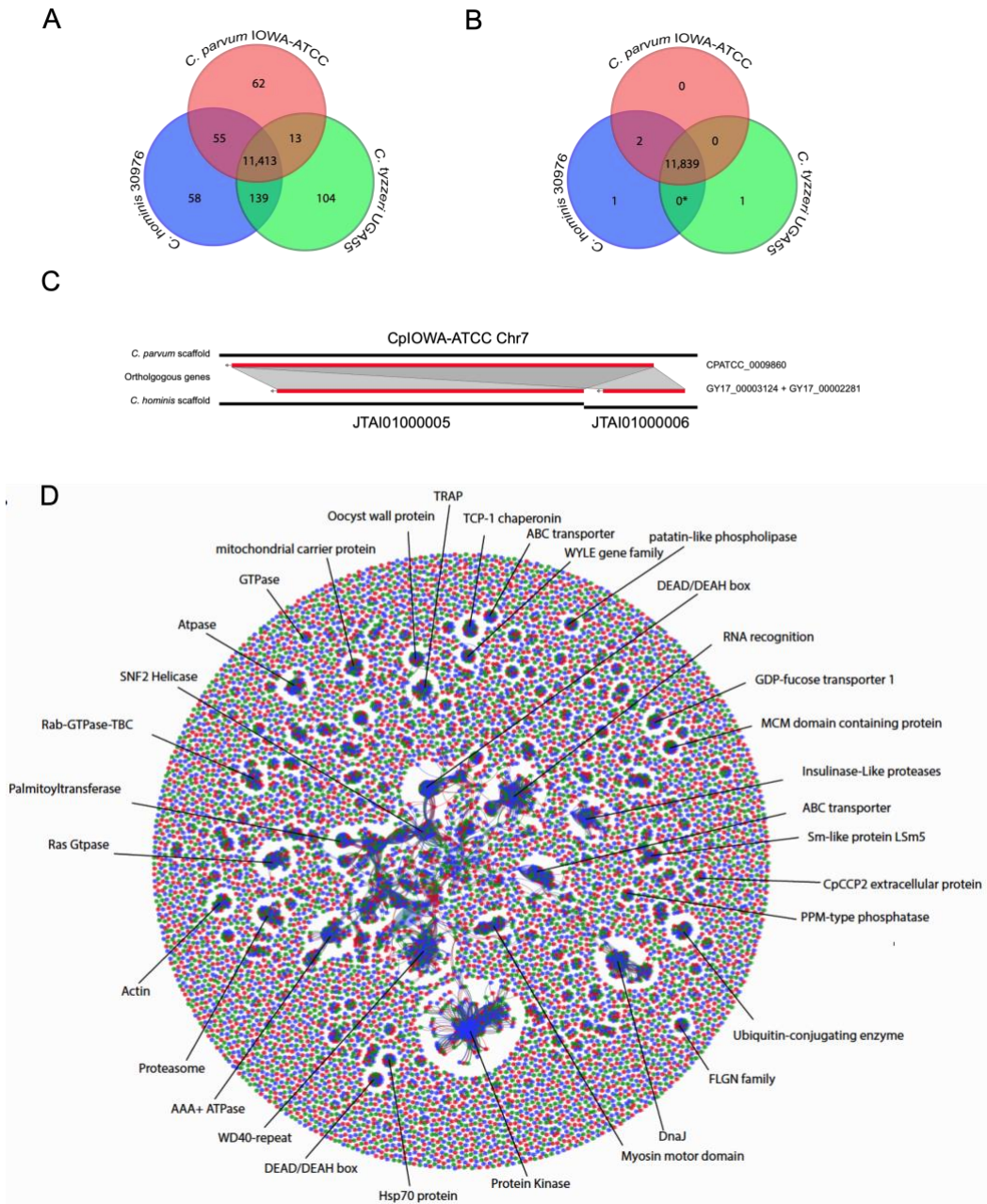
231

**Figure 3.** Comparative genome-wide synteny between six different species of *Cryptosporidium*. Highly conserved regions between the genomes are colored in order from red (5' end of chromosome 1) to blue (3' end of chromosome 8) with respect to genomic position of *C. parvum*. The cladogram topology was determined via a maximum likelihood analysis of 2700 revisited single copy orthologs. Animal icons represent the major hosts for these species. \**C. parvum* and *C. ubiquitum* are zoonotic with many hosts.

232

233 The consistent annotation of the species closest to *C. parvum* IOWA-ATCC, *C. hominis*  
234 30976 and *C. tyzzeri* UGA55, permitted the detection of differences in protein encoding gene  
235 content and copy number variation. An automated orthology analysis between all three gene sets  
236 revealed that, ~94% of the genes were conserved among all species. Of the 4,008 ortholog  
237 groups identified, most annotated gene families were maintained with a similar number of  
238 paralogs (max = 6) detected in the same ortholog group, but the number of singletons varied  
239 between the three species (Fig. 4A; Supplemental Table S7). Some of these post-comparative  
240 annotation gene differences appeared to be unique to a particular species (Supplemental Table  
241 S8). Of the 224 singletons detected, we observed only 0, 1 and 1 potential truly species-specific  
242 genes in *C. parvum* IOWA-ATCC, *C. hominis* 30976 and *C. tyzzeri* UGA55, respectively following  
243 manual inspection (Fig. 4B). Both species-specific genes are uncharacterized proteins. The  
244 remaining 253 singletons are detected but incomplete in the fragmented assemblies of *C. hominis*  
245 and *C. tyzzeri*, appearing as split genes, frame-shifts, missed calls near a gap or contig break and  
246 putative false gene predictions in small contigs (Fig. 4C; Supplemental Fig. S3). The majority of  
247 gene content differences between these species are gene copy number variations and not gene  
248 presence or absence.

249



250

**Figure 4.** Ortholog distribution of protein-encoding genes. (A) Venn diagram of orthologous gene sequences between three closely related *Cryptosporidium* species (pre-investigation); (B) Venn diagram of same orthologous gene sequences (post-investigation) following removal of false species-specific genes, e.g. artifacts. \*The 139 genes shared between *C. hominis* and *C. tyzzeri* in panel A are in complex regions with repeats and gaps and do not have enough evidence to prove their uniqueness at this stage given the available assemblies, so they are considered artifacts at this time; (C) Orthology based synteny overview of Chr7 singleton-like, putative paralog artifact generated by a split gene due to the genome assembly fragmentation in one species; and (D) Graphical representation of ortholog clusters between the three closely related



*Cryptosporidium* species. *C. parvum* IOWA-ATCC: red; *C. hominis* 30976: blue; *C. tyzzeri* UGA55: green.

251 We mapped Illumina reads from *C. parvum* IOWA, *C. hominis* TU502-2012 and *C. tyzzeri*  
252 UGA55 to the new *C. parvum* IOWA-ATCC long-read assembly to identify and assess putatively  
253 overly collapsed regions (repetitive regions represented by only a single repeat in the assembly)  
254 (Supplemental Table S1; Supplemental Fig. S4). Our pipeline detected 14 compressions > 100  
255 bp in length in the *C. parvum* IOWA II genome assembly compared to 8 in the new *C. parvum*  
256 IOWA-ATCC assembly. These compressions are not always related to genic regions and vary in  
257 genome location and predicted copy number. Some of these apparently collapsed regions, were  
258 conserved between both *C. parvum* assemblies but varied in different species (Supplemental Fig.  
259 S5). The collapsed genic regions are composed of rRNA genes, some uncharacterized proteins,  
260 GMP synthase, aspartate-ammonia ligase, tryptophan synthase beta and MEDLE genes. Most of  
261 the observed and fixed compressions do not contain any annotated genes.

262

263

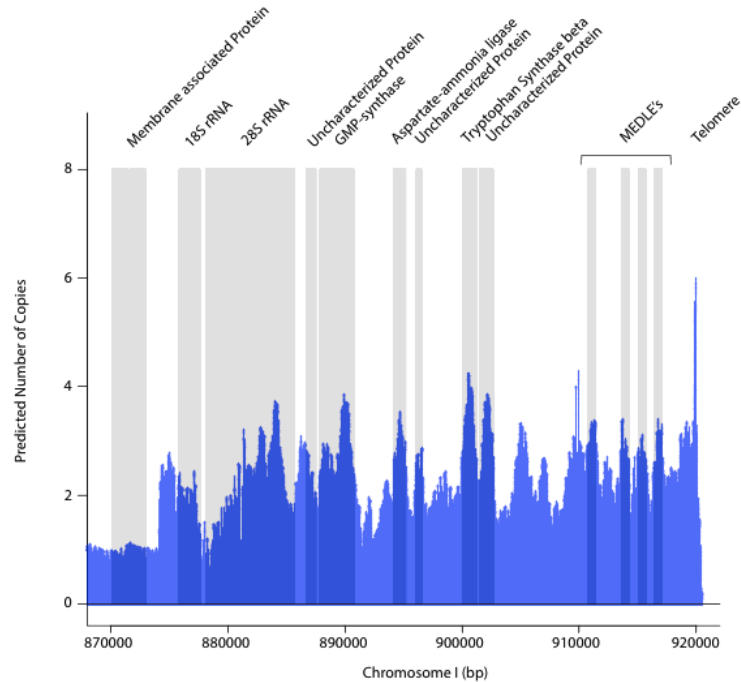
### 264 **A closer look at subtelomeric regions reveals their complexity and relevant biology**

265

266 As shown in the read depth coverage analysis and in Supplemental Table S1, the new  
267 assembly was able to fix most of the collapsed regions in the *C. parvum* IOWA-ATCC genome.  
268 Interestingly, one subtelomeric region in Chr1 still has compressions suggesting that most of the  
269 genes present in this region have more than one copy (Fig. 5). This region reveals at least 13  
270 genes which vary in copy number between different *Cryptosporidium* species (Supplemental Fig.  
271 S5). The genes contained in this region are 18S rRNA, 5S rRNA and 28S rRNA, uncharacterized  
272 proteins, a GMP synthase, an aspartate-ammonia ligase, tryptophan synthase beta and a cluster  
273 of several MEDLE genes. Some of these genes, such as the tryptophan synthase beta and the  
274 MEDLE's are the focus of considerable research since they may be related to parasite survival  
275 and are potentially involved in parasite invasion, respectively (Sateriale and Striepen 2016; Li et  
276 al. 2017; Fei et al. 2018). The number of copies predicted here for the rRNAs and MEDLE's are  
277 underrepresented as they also have paralogs on Chr 2 and Chr 5, respectively.

278

279



280

**Figure 5.** Chromosome 1 subtelomeric region read depth coverage plot normalized by single copy genes. Illumina reads from *C. parvum* IOWA-ATCC DNA are mapped to the *C. parvum* IOWA-ATCC long-read assembly to identify read pileups and estimates of sequence copy number. Vertical grey areas indicate regions with annotated genes. The GMP-synthase shaded region also contains a small uncharacterized protein.

281

282

283 Since we have an apparent compression in a subtelomeric region assembly with no gaps

284 and good PacBio long read coverage, we hypothesized that these extra copies might derive from

285 unassembled regions. The chromosomal-level IOWA-ATCC assembly was only missing three

286 telomeric regions, both ends of Chr 7 and one telomere of Chr 8. Using existing PacBio long-

287 reads we were able to identify a few reads that extended into rRNA regions on the chromosomes

288 missing telomeres. We attempted re-assembly with only PacBio reads and we could not

289 convincingly resolve the missing regions. Thus, we generated very deep (1200 X) Oxford

290 Nanopore (ONT) single molecule reads from *C. parvum* IOWA-BGF (ATCC was not available).

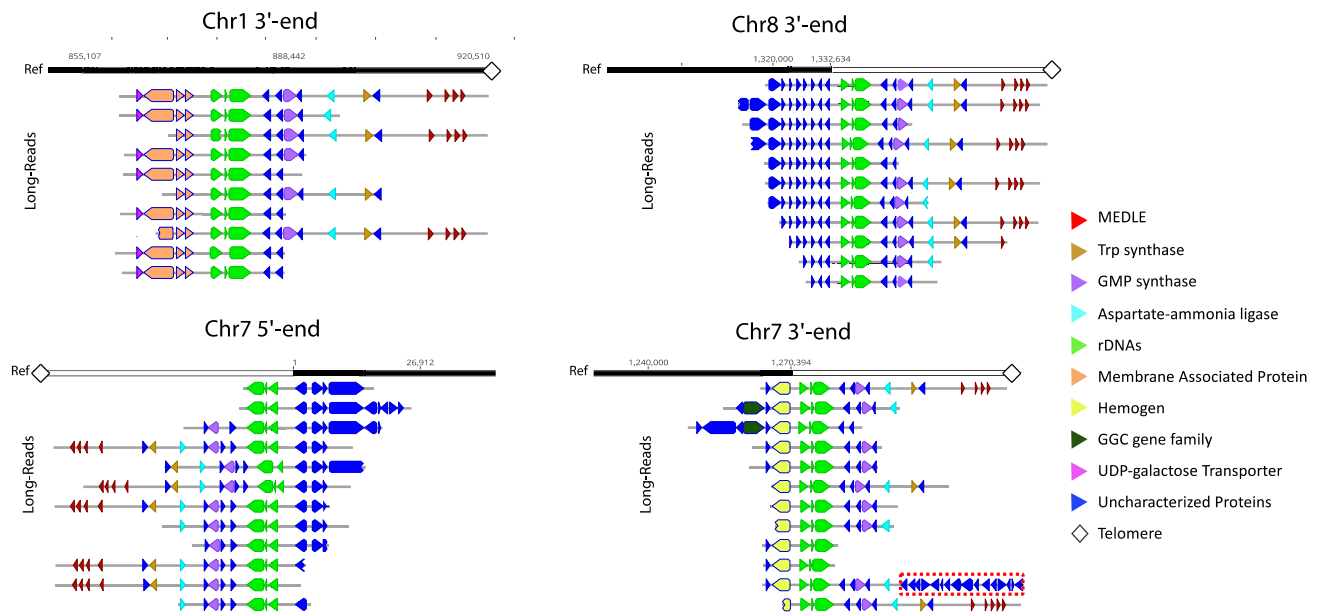
291 The ONT reads revealed related, yet unique subtelomeric regions linked to the chromosomes

292 missing their telomeres, in addition to Chr 1 (Fig. 6). We found good ONT long-read support for

293 these regions. Notably, each different subtelomeric region is flanked by ribosomal RNAs and we

294 also note that there is slight variation observed among the reads for each chromosome end.

295



296

**Figure 6.** Related subtelomeric regions on different *C. parvum* chromosomes are supported by ONT long-reads. Individual ONT long-reads provide evidence of at least four different, yet related, subtelomeric regions that extend into the chromosomes that were missing telomeres (Chr 7 and Chr 8) in addition to Chr 1. The white and black reference bar above each collection of annotated Nanopore reads identify the newly identified subtelomeric regions (white) and existing assembly (black). The red box on the penultimate read on the Chr 7 3' end panel indicates a unique region of insertion (nucleotide positions 1191705-1217462). This region contains mostly uncharacterized proteins and two transferases. Each ONT read is annotated as indicated in the key.

297

298

299

### Fast evolving genes in *C. parvum*

300 The new gapless genome assembly and annotation presented an opportunity to revisit the  
 301 prediction of fast-evolving genes in this species. We performed a Single Nucleotide Variant, SNV,  
 302 analysis using 136 different *C. parvum* WGS data sets obtained from GenBank (Supplemental  
 303 Table S9) using the new assembly and annotation. A total of 24,407 positions were found to  
 304 contain at least one high-confidence bi-allelic variant. Multiallelic calls were removed to guard  
 305 against mixed infections. The biallelic variants reflect 3892 genes, 190 of which show a Ka/Ks  
 306 ratio of non-synonymous/synonymous rates of > 1.0 (Supplemental Table S10). Of the 190, 24  
 307 genes were previously identified and 124 are classified as uncharacterized proteins, 93 of which  
 308 are annotated as having a signal peptide or being secreted. All previously identified fast evolving  
 309 genes were detected, including: Insulinase-like protein (CPATCC\_0017080), an uncharacterized  
 310 secreted protein (CPATCC\_0010380), *gp60* (CPATCC\_0012540) and others (Strong et al. 2000;  
 311 Sanderson et al. 2008; Nader et al. 2019; Zhang et al. 2019). The top eight genes by Ka/Ks ratio  
 312 have not been previously reported. Gene family members such as MEDLEs, FLGN and SKSR  
 313 were also detected but significantly, new members of each of these families are identified as also  
 314 under positive selection. A family of WYLE (Sanderson et al. 2008) proteins is also identified as  
 315 under selection.

316

317

### DISCUSSION

318  
319 The first genome assembly sequence of *Cryptosporidium parvum* IOWA II (Abrahamsen  
320 et al. 2004) was excellent given the technology at the time and because of its quality the  
321 community has relied on this genome assembly and annotation to design their experiments.  
322 However, gaps and ambiguous bases remained, and there was little available expression or  
323 orthology evidence to assist the annotation. We used PacBio, Nanopore and Illumina sequencing  
324 technologies to generate a new complete genome assembly of *C. parvum* strain IOWA-ATCC.  
325 We then applied de novo and evidence-based annotation approaches with manual curation of two  
326 additional species to generate consistent annotation that could be used to detect unique genes  
327 and genomic differences between species and strains.

328 The first, expected, finding was that the *C. parvum* IOWA stain is continuing to evolve  
329 (Cama et al. 2006) as it is maintained by passage through cattle in a few different locations for  
330 research use. Some natural *Cryptosporidium* isolates have been propagated in unnatural hosts  
331 before sequencing. Thus, selection during propagation or maintenance via animal propagation  
332 may lead differences relative to circulating parasites. This phenomenon has been observed in  
333 other protozoan parasites (Lecomte et al. 1992; Sutherland et al. 1996; Akiyoshi et al. 2002; Chan  
334 et al. 2015; Isaza et al. 2015). Genomic DNA for the 2004 *C. parvum* IOWA II and *C. parvum*  
335 IOWA-ATCC were obtained from the same source, but many years apart. We note small  
336 differences in the *gp60* sequence, and an overall genome average difference of ~0.07% in identity  
337 (Supplemental Table S2). Changes were also observed in *Plasmodium* species, after being  
338 propagated for a long time period (Claessens et al. 2017) and is associated with loss of infectivity  
339 and virulence in some strains (Segovia et al. 1992).

340 When compared to *C. hominis* and *C. tyzzeri*, which are 95–97% identical in available  
341 nucleotide sequence, incongruences in the annotated gene models with respect to the new *C.*  
342 *parvum* IOWA-ATCC genome assembly were obvious. The differences result in part from some  
343 genome assemblies that contain numerous sequence gaps and little experimental evidence (i.e.  
344 RNA-Seq data for all major developmental stages) to permit accurate annotation. The gaps can  
345 lead to pseudo-gene annotations or split genes due to frame-shift artifacts. For example, regions  
346 with gaps or errors in the base call can lead to false stop codons, or frameshifts that are usually  
347 detected as incomplete pseudogenes or a gap can cause a predicted gene to be split into more  
348 than one piece. Additionally, *in silico* prediction tools are usually not trained for non-model  
349 organisms and *C. parvum* is so distant from other sequenced organisms, there is little synteny or  
350 orthology to help guide the various efforts. These mis-annotations can be detected and avoided  
351 if an evidence and homology-based curation between different samples is conducted.

352 As observed in Table 2, annotated protein-coding gene numbers do not exactly match  
353 between the closely related species with high sequence identity. This difference is explained by  
354 the gaps in the comparator genome assemblies for *C. hominis* and *C. tyzzeri*. These gaps  
355 interrupt the open reading frames (ORFs) causing split genes and frame shifts. Thus, some gene  
356 models are not necessary missing in an organism. These differences affect similarity-based  
357 analyses such as ortholog detection, giving the wrong impression that some of these partially  
358 annotated genes are unique for a species (Fig. 4; Supplemental Table S8). These mis-  
359 interpretations can sabotage some experimental designs that may use an incorrect basis for  
360 experimental design or analysis (Baptista and Kissinger 2019). These regions with problems are  
361 usually complex and some have high polymorphism rates (e.g., positive selection). So, false  
362 assumptions regarding species-specific genes can affect many downstream analyses including  
363 the detection of highly polymorphic loci.

364 In this study we were able to improve the structural and functional annotation of these  
365 genome assemblies, by using two different approaches: (i) using seven full-length stranded cDNA  
366 libraries derived from three time points (0h, 24h and 48h post infection) increasing the expression  
367 percentage representation of *C. parvum* and transcriptome data from RNA-seq analyses were  
368 generated to improve the gene models deposited in CryptoDB.org (Tandel et al. 2019); and (ii)

369 by using homology information to construct a consistent genome annotation between three  
370 different close-related species. This approach facilitated a proper comparative analysis of genome  
371 content differences between the species compared. Our analyses reveal that the compared  
372 species only differ slightly in gene content for the regions that can be compared. Most differences  
373 are related to slight structural variation, such as small translocations and inversions, and by copy  
374 number variation as revealed by read depth coverage analysis. Previous studies have reported a  
375 lack of DNA methylation in *Cryptosporidium* and other parasites (Gissot et al. 2008). The *C.*  
376 *parvum* C-5 cytosine-specific DNA methylase (*Dnmt2*) sequence was previously annotated as  
377 truncated (Abrahamsen et al. 2004; Isaza et al. 2015) and lacking a DNMT-specific motif  
378 containing a prolyl-cysteiny dipeptide (Abrahamsen et al. 2004; Ponts et al. 2013; Isaza et al.  
379 2015). The new *Cryptosporidium parvum* IOWA-ATCC whole genome assembly and annotation  
380 reveals a complete ortholog of the *Dnmt2* DNA methylase family. The lack of this N-terminus has  
381 been cited as a possible reason for the lack of DNA methylation in *C. parvum* (Ponts et al. 2013).

382 Apicomplexans have reductive streamlined genomes, that range from ~8.5 to ~125  
383 megabases. *Cryptosporidium* species have among the most compacted genomes, with 504 bp  
384 average length between the stop codon of one gene and the start codon of the next gene.  
385 *Cryptosporidium* also has few protein-encoding genes (~3950) relative to other apicomplexans  
386 with up to ~8000 (Kissinger and DeBarry 2011). Studies shows that *Cryptosporidium* may have  
387 adapted a novel type of nucleotide transporter for ATP uptake from the host (Pawlowic et al. 2019).  
388 Given the compactness of this parasite genome sequence, the gene loss may be compensated  
389 for by the higher number of transporters found in our re-analysis. These findings will facilitate  
390 future studies of alternative metabolic pathways to better understand the biology and evolution of  
391 parasitism of this organism.

392 Chromosomal inversions are known to affect rates of adaptation, speciation, and the  
393 evolution of chromosomes (Guo et al. 2015). Comparative genomic studies and population  
394 models for several organisms, suggests that inversions can spread by suppressing recombination  
395 between loci and generating areas of linkage disequilibrium. Local adaptation mechanisms  
396 applied to demographic and genetic situations, can drive inversion to high frequency if there is no  
397 countervailing force, thus explaining fixed differences observed between populations and species  
398 (Kirkpatrick and Barton 2006). Previous studies identified potential chromosomal inversion sites  
399 between *Cryptosporidium* species relative to *C. parvum* IOWA II (Guo et al. 2015; Isaza et al.  
400 2015). The new long-read genome assembly of *C. parvum* IOWA-ATCC revealed some potential  
401 inversion sites, in chr 2, chr 4 and chr 5, that are flanked by poorly sequenced and gapped regions  
402 in some species, (Piper et al. 1998; Bankier et al. 2003). Since the other species still lack physical  
403 evidence for their chromosomal structures, further long-read sequencing or chromosome  
404 conformation capture sequencing, such as Hi-C, is still needed to detect and validate species-  
405 specific structural variations for the other *Cryptosporidium* species.

406 The lack of three telomeres in the new high-quality long-read assembly was an intriguing  
407 result that can be explained by the detection of three putative similar but not identical copies of  
408 subtelomeric regions containing genes including tryptophan synthase beta, the MEDLE genes  
409 and 18S/28SrRNA cluster among others. This finding raises the possibility of this species having  
410 misincorporation of telomers by its telomerase, as was observed in other protists (McCormick-  
411 Graham et al. 1997) or recombination between telomeres by break-induced replication, such as  
412 has been observed in yeasts (McEachern and Iyer 2001; McEachern and Haber 2006), and  
413 telomere maintenance by recombination as is observed in human cancers (Natarajan et al. 2006).  
414 Since some of the genes in this region are possibly essential genes for parasite survival (Sateriale  
415 and Striepen 2016), the fact that they may exist in multiple copies and can possibly generate  
416 variation as a result of recombination could explain an alternate new survival mechanism in this  
417 streamlined parasite genome. We have support from single molecule sequencing that indeed this  
418 region is detected on 4 different chromosome ends (Fig. 6). This potential subtelomeric plasticity  
419 resulting in a possible transfer of important gene sequences between homologous and

420 nonhomologous chromosome ends, could affect genetic manipulations and may affect phenotype.  
421 We believe that these structures are varying within the *Cryptosporidium* population, which is hard  
422 to detect, since we do not yet have any evidence that all 4 related chromosome ends are present  
423 in a single cell. Thus, the Nanopore reads may be representing population level variation, which  
424 also raises the possibility of recombination or gene conversion as *Cryptosporidium* requires  
425 sexual recombination to form excreted oocysts. Currently, cloning does not exist for  
426 *Cryptosporidium*. Thus, oocysts used for sequencing must be considered a population even if  
427 sequence is derived from single cell sequencing (Troell et al. 2016) as oocysts still contain four  
428 haploid meiotic progeny (sporozoites). A truly single-cell approach, which will facilitate  
429 recombination and sub-telomeric plasticity studies, will require single-sporozoite sequencing, but  
430 this is still impossible in the absence of genome amplification.

431 *Cryptosporidium* species are usually typed and characterized by the community using a  
432 small number of genetic markers including 18S, COWP, HSP70, and *gp60* (Ghaffari et al. 2014).  
433 As shown in this study *gp60* which is a fast-evolving gene used for *Cryptosporidium* subtyping  
434 characterization, had small differences between *C. parvum* IOWA II and *C. parvum* IOWA-ATCC.  
435 The parasites used to generate these sequences originated from the same propagated strain but  
436 were collected at different times. Using just one marker to characterize an obligately sexual  
437 organism with 8 chromosomes is problematic. In this study, we confirm an existing group of fast  
438 evolving genes and identify 166 additional potential candidates distributed across all 8  
439 chromosomes. Some of these genes belong to gene families so to avoid artifacts only uniquely  
440 mapped reads were used for the SNV analysis. The genes identified here can be used to help  
441 the community develop additional markers with better resolution for typing parasite isolates. Given  
442 that only 136 isolates from a small geographic region have been sampled, the potential to identify  
443 additional genes is high. Newer techniques such as hybrid capture bait set techniques  
444 (Mamanova et al. 2010) are a powerful future alternative to characterize and select  
445 *Cryptosporidium* population variants and better characterize genetic diversity.

446  
447 The new *C. parvum* long-read assembly combined with a consistent comparative  
448 annotation has proven incredibly powerful. The species analyzed here have different host  
449 preferences and pathogenicity. Comparisons of previous sequences and annotation suggested  
450 numerous gene content differences. However, this systematic study reveals that the primary  
451 differences between the zoonotic *C. parvum*, the anthroponotic *C. hominis* and the rodent-  
452 infecting *C. tyzzeri* are SNVs and CNVs rather than differences in unique gene content. Finally,  
453 new findings related to within parasite and/or within population subtelomeric amplification and  
454 variation events in *C. parvum* reveal a new level of genome plasticity that will impact some genetic  
455 manipulations and may affect the organisms' phenotype.

456  
457

## 458 **METHODS**

459

### 460 **Sample DNA source and dataset used.**

461

462 *Cryptosporidium parvum* IOWA-ATCC DNA from oocysts/sporozoites was purchased from  
463 the ATCC. The source was the University of Arizona, Sterling Parasitology Laboratory. Its GP60  
464 subtype (IIa) is the same as the current *C. parvum* IOWA II reference genome sequence also  
465 used in this work. *Cryptosporidium parvum* DNA was also prepared from oocysts obtained in 2018  
466 from Bunch Grass Farms, Deary, ID. This isolate is referred to as IOWA-BGF in this study. The  
467 *C. hominis* 30976 and UdeA01 genome assemblies, are human isolates. The *C. tyzzeri* assembly  
468 a natural mouse model of Cryptosporidiosis. The 136 *C. parvum* sample accession numbers used  
469 for the positive selection analysis are available in Supplemental Table S9.

470

## 471 ***Cryptosporidium parvum* IOWA-ATCC sequencing and genome assembly**

472  
473 PacBio RSII and Illumina HiSeq 2000 sequencing were both performed at the Wellcome  
474 Sanger Institute, UK. The *Cryptosporidium parvum* IOWA-ATCC reads were first assembled using  
475 the PacBio open source SMRTlink v6.0 from 9 PacBio SMRT cells, with ~75x mean genome  
476 coverage. The resulting assembly was then submitted to the accuracy improver tool Sprai  
477 0.9.9.23 (<https://sprai-doc.readthedocs.io/en/latest/index.html>) and then had gaps filled using  
478 PBJelly 15.24.8 (English et al. 2014) using PacBio reads and IMAGE 2.4.1 (Swain et al. 2012)  
479 with Illumina reads. A manual inspection and improvement using GAP5 (Bonfield and Whitwham  
480 2010) was needed to better access complex regions, and the final scaffolded genome assembly  
481 was polished with Illumina reads using iCORN2 0.95 (Otto et al. 2010) and Pilon 1.22 (Walker et  
482 al. 2014).

483 Oxford Nanopore (ONT) single molecule long-read sequencing was performed on DNA  
484 from *C. parvum* IOWA-BGF (The ATCC®PRA-67DQ™ ran out of stock) following the protocol  
485 recommended by for an R9.4.1 flow cell. MinION ONT sequencing was performed at the Georgia  
486 Genomics Bioinformatics Core (GGBC) at the University of Georgia, USA, using an R.9.4 flow  
487 cell and the rapid sequencing kit (SKT-RAD004). The ONT long-reads generated >1000x  
488 coverage of the *Cryptosporidium parvum* genome. This high coverage complemented the PacBio  
489 data to confirm and fix several complex regions. The final assembly was submitted with the current  
490 reference and genome assemblies of other closely related species to QUAST v.5.02 (Gurevich et  
491 al. 2013) to compare and evaluate the quality of the new genome assembly.

492

## 493 ***Cryptosporidium* genome reannotation**

494 Genome annotation was generated with: (a) an ab initio prediction using GeneMark-ES  
495 4.57 (Lomsadze et al. 2005); (b) evidence-trained predictions by SNAP/Maker (Cantarel et al.  
496 2008; Johnson et al. 2008) and (c) Augustus (Stanke and Morgenstern 2005). For training, we  
497 used publicly available data from each respective species: RNA-seq (strand and non-strand  
498 specific), ESTs, previously predicted proteins and MassSpec proteomics data when available. In  
499 parallel we also generated transcriptome assemblies using HISAT2 v.2.1.0 (Kim et al. 2015) and  
500 StringTie v.1.3.4 (Pertea et al. 2015), and non-coding RNA predictions were generated for *C.*  
501 *parvum* as described (Li et al. 2020). Manual curation of all genes in the context of existing  
502 molecular evidence was performed using a local installation of WebApollo2 (Lee et al. 2013).

503 As each genome species analyzed has a different number of publicly available data sets,  
504 we also used each curated genome annotation in comparison with the others using the Artemis  
505 Comparison tool (ACT) 17.0.1 (Carver et al. 2005), allowing us to perform comparative annotation  
506 and resolve discrepancies via homology. All protein-encoding genes annotated for each genome  
507 sequence were submitted to OrthoFinder v.2.3.7 (Emms and Kelly 2015) to detect paralogs,  
508 orthologs and singletons. All singletons were then selected for a comparative manual curation  
509 using MCScanX 0.8 (Wang et al. 2012) and JBrowse (Buels et al. 2016) between all three species  
510 to verify their uniqueness and assess the contribution of sequence gaps or misassembly to the  
511 findings. We considered the following error types: Split genes caused by frameshifts or early stop-  
512 codons, lack of stranded RNAseq to confirm the gene model, and the presence of a gapped region  
513 in the genome assembly. All genes that did not fall into one of these categories were considered  
514 to be unique.

515

## 516 **Functional annotation**

517

518 Following structural annotation, the predicted protein sequences were used to search  
519 Swiss-pro curated (sprot) and not-curated (Trembl) and the NCBI non-redundant Protein  
520 database with BLASTP and an e-value threshold at the superfamily level of 1e-6. Protein  
521 structure similarity was explored using I-TASSER (Roy et al. 2010). Protein sequences were

522 divided into two major groups distributed according their length for the I-TASSER analysis: (i)  
523 peptide sequences < 750 aa; and (ii) shorter sequential segmental peptide sequences < 750 aa  
524 derived from annotated proteins > 750aa. Structures were predicted for each peptide using the I-  
525 TASSER suite and aligned to solved crystal structures in the protein data bank (PDB) using the  
526 cofactor algorithm (Roy et al. 2012). InterPro codes were assigned to the query peptide  
527 sequence via InterProScan v.5.23-62.0 (Quevillon et al. 2005) using 11 different default  
528 databases. The PFAM codes available for PDB crystal structures were transposed to InterPro  
529 codes using the R pfam.db library. The presence of at least one matching InterPro code  
530 assigned to both the query and the reference peptides was taken to indicate a greater likelihood  
531 of structural similarity of the predicted structure and considered “high-confidence”. A random  
532 forest classifier was trained to distinguish between a test set of high- and low-confidence  
533 models, and was then applied to the entire predicted proteome to identify additional high-  
534 confidence-like models among unannotated proteins, as described in Ansell et. al 2019 (Ansell  
535 et al. 2019). BLAST2GO (Conesa et al. 2005) version 4.1.9 was used to assign Enzyme Code  
536 (E.C) and Gene Ontology (GO) terms. Following this functional annotation, we compared the  
537 existing protein product names to the new functional results. Some structural information, such  
538 as protein domain and repeat pattern content were added to some uncharacterized proteins and  
539 nomenclature errors were corrected according to the NCBI annotation submission guide.

540

#### 541 **Transporter prediction**

542

543 Predicted proteins were submitted to four different transporter prediction methods: (i) local  
544 alignment using BLASTP against TCDB (Saier et al. 2009) transporter proteins with a threshold  
545 e-value of 1e-5 cutoff to find potential transporter similarities; (ii) TMHMM (Server v. 2.0) (Krogh  
546 et al. 2001) and SignalP (Server 4.1) (Bendtsen et al. 2004) was applied to reduce false positives  
547 from the TCDB blast results. Transporter candidates with no transmembrane domains or  
548 candidates with only one transmembrane prediction while having signal peptides predicted were  
549 removed; (iii) TransAAP (Ren et al. 2007), which is a TC-based (Transporter Classification from  
550 TCDB) transporter annotation tool on the TransportDB v2.0 website (Ren et al. 2007), that was  
551 used to provide information about potential transporter identity and substrate; and (iv) a structural  
552 proof for candidate transporters using Phyre2.0 (Kelley et al. 2015). Final candidate transporters  
553 were checked according to above results as well as annotations obtained from InterProScan 5.44  
554 (Jones et al. 2014).

555

#### 556 **Comparative and phylogenetic analysis**

557

558 Comparative genome-wide synteny between *Cryptosporidium* species was performed  
559 using Murasaki v.1.68.6 (Popendorf et al. 2010) with default settings. The cladogram topology  
560 was determined via a maximum likelihood analysis of 2700 single copy orthologs using JTT+I as  
561 the substitution model as predicted by Modeltest-NG (Darriba et al. 2020). The consensus tree  
562 was constructed from 1000 bootstrap replicates. The consistency of annotation and potential  
563 gene family copy number variations (CNVs), were determined with Orthofinder v.2.2.7 (Emms  
564 and Kelly 2015) which identified all orthologs and paralogs. Orthofinder BLASTP results were  
565 parsed to examine the relationships between proteins using an e-value threshold of 1e-20 and  
566 identities > 35% between protein pairs longer than 100 amino-acids. The data were visualized  
567 using Gephi (<https://gephi.org/>) with the Fruchterman-Reingold layout.

568 Copy number variation was also determined by aligning Illumina sequence reads from  
569 each closely related species studied to the new *C. parvum* IOWA-ATCC reference genome  
570 sequence to check for potential CNV regions by looking for variations in read depth coverage.  
571 The alignment was performed using BWA mem 0.7.17 (Li and Durbin 2009) with default options



572 and the alignment depth per base was calculated using BEDTools genomecov 2.29.2 (Quinlan  
573 and Hall 2010) and SAMtools depth 1.6 (Li et al. 2009).

574

### 575 **Resolving the structure of repetitive subtelomeric regions**

576

577 Following the CNV analysis, the sequence content of the putatively compressed regions  
578 and their non-compressed sequence boundaries of the *C. parvum* IOWA-ATCC assembly were  
579 used to build a BLAST database. We then selected single oxford nanopore single molecule reads  
580 using BLASTn 2.10.0 (Camacho et al. 2009) to detect sequences capable of aligning to  
581 compressed regions and then determine their putative assembly structures. Following ONT read  
582 selection, the ONT reads were polished with Illumina reads using proovread 2.14.1 (Hackl et al.  
583 2014) and Pilon 1.22 (Walker et al. 2014). To map these polished reads against the genome  
584 assembly and avoid bias/competition between sites, all putatively compressed genome assembly  
585 regions were artificially split into fragments, effectively making the chromosomes with compressed  
586 regions fragmented. Reads were aligned to all chromosome fragments using the Geneious  
587 mapper 2019.1.3 (<https://www.geneious.com>) with medium-sensitivity and those chromosome  
588 fragments with hits were annotated and analyzed for validation and verification of their structure.

589

### 590 **Variant analysis, selection prediction and populational analysis**

591

592 Illumina sequence reads from 136 different isolates of *C. parvum* from different  
593 geographical locations (Supplemental Table S9) were aligned against the *C. parvum* IOWA-ATCC  
594 reference genome sequence using BWA-MEM (Li and Durbin 2009), the bam files were parsed  
595 to select uniquely mapped reads and to mark duplicates and remove redundancy using PICARD  
596 (Broad\_Institute) and then submitted to a Variant call analysis using GATK 3.8 Haplotypecaller  
597 (McKenna et al. 2010). These results were then filtered by mapping quality > 40 and depth  
598 coverage >10. Because mixed infections exist, we restricted analysis to biallelic sites. The  
599 individual VCF files were combined into one GVCF file using the GATK tool GenotypeGVCF. After  
600 selecting just single nucleotide variants (SNVs) from this data, the combined gvcf file was  
601 annotated using the software snpEff v.4.3 (Cingolani et al. 2012). The number of synonymous  
602 and non-synonymous variants were taken from the annotated gvcf file and parsed to calculate  
603 the Ka/Ks ratio of non-synonymous/synonymous rates. Genes with ratios > 1.5, indicative of  
604 positive selection, were detected and denoted as fast evolving genes within the *C. parvum*  
605 population.

606

607

608

609

### 610 **DATA ACCESS**

611

612 The sequencing data, genomes and annotation generated in this study have been submitted to  
613 the NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession  
614 numbers PRJNA573722, PRJNA252787, PRJEB3213 and PRJNA388495. *C. hominis* UdeA01  
615 assembly and TU502 Illumina reads used are in BioProjects PRJEB10000 and PRJNA222836,  
616 respectively. The data are also available at CryptoDB.org (Heiges et al. 2006).

617

### 618 **ACKNOWLEDGMENTS**

619

620 This work was supported by Bill and Melinda Gates Foundation grant OPP1151701 to JCK, The  
621 Wellcome Trust via its core funding of the Wellcome Sanger Institute (grant WT206194) and  
622 NHMRC Investigator Grant (APP1194330) to ARJ.

623

## 624 **AUTHORS CONTRIBUTIONS**

625

626 RPB and JCK designed research; RPB and JCK performed research; AS, JD and BS contributed  
627 with new reagents and samples; BA and AJ contributed with analytical tools; MS, KB, AT, MB and  
628 JAC contributed Illumina and PacBio sequencing; RPB, YL, KB, AT, RX, EDS, GWC and JCK  
629 analyzed data; RPB and JCK wrote the paper and ARJ, BREA, BS, AS and JAC provided  
630 feedback.

631

## 632 **DISCLOSURE DECLARATION**

633

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

635

## 636 **BIBLIOGRAPHY**

637

638 Abrahamsen MS, Templeton TJ, Enomoto S, Abrahante JE, Zhu G, Lancto CA, Deng M, Liu C,  
639 Widmer G, Tzipori S et al. 2004. Complete genome sequence of the apicomplexan,  
640 *Cryptosporidium parvum*. *Science* **304**: 441-445.

641 Akiyoshi DE, Feng X, Buckholt MA, Widmer G, Tzipori S. 2002. Genetic analysis of a  
642 *Cryptosporidium parvum* human genotype 1 isolate passaged through different host  
643 species. *Infect Immun* **70**: 5670-5675.

644 Ansell BRE, Pope BJ, Georgeson P, Emery-Corbin SJ, Jex AR. 2019. Annotation of the *Giardia*  
645 proteome through structure-based homology and machine learning. *Gigascience* **8**.

646 Bankier AT, Spriggs HF, Fartmann B, Konfortov BA, Madera M, Vogel C, Teichmann SA, Ivens  
647 A, Dear PH. 2003. Integrated mapping, chromosomal sequencing and sequence  
648 analysis of *Cryptosporidium parvum*. *Genome Res* **13**: 1787-1799.

649 Baptista RP, Kissinger JC. 2019. Is reliance on an inaccurate genome sequence sabotaging  
650 your experiments? *PLoS Pathog* **15**: e1007901.

651 Bendtsen JD, Nielsen H, von Heijne G, Brunak S. 2004. Improved prediction of signal peptides:  
652 SignalP 3.0. *Journal of molecular biology* **340**: 783-795.

653 Bonfield JK, Whitwham A. 2010. Gap5--editing the billion fragment sequence assembly.  
654 *Bioinformatics (Oxford, England)* **26**: 1699-1703.

655 Broad\_Institute. Picard Tools.

656 Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein DM, Elsik CG, Lewis  
657 SE, Stein L et al. 2016. JBrowse: a dynamic web platform for genome visualization and  
658 analysis. *Genome Biol* **17**: 66.

659 Cama VA, Arrowood MJ, Ortega YR, Xiao L. 2006. Molecular Characterization of the  
660 *Cryptosporidium parvum* IOWA Isolate Kept in Different Laboratories. *J Eukaryot*  
661 *Microbiol* **53 Suppl 1**: S40-42.

662 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.  
663 BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 421.

664 Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, Holt C, Sanchez Alvarado A, Yandell  
665 M. 2008. MAKER: an easy-to-use annotation pipeline designed for emerging model  
666 organism genomes. *Genome Res* **18**: 188-196.

667 Carver TJ, Rutherford KM, Berriman M, Rajandream MA, Barrell BG, Parkhill J. 2005. ACT: the  
668 Artemis Comparison Tool. *Bioinformatics (Oxford, England)* **21**: 3422-3423.

- 669 Chalmers RM, Smith R, Elwin K, Clifton-Hadley FA, Giles M. 2011. Epidemiology of  
670 anthroponotic and zoonotic human cryptosporidiosis in England and Wales, 2004-2006.  
671 *Epidemiology and Infection* **139**: 700-712.
- 672 Chan ER, Barnwell JW, Zimmerman PA, Serre D. 2015. Comparative analysis of field-isolate  
673 and monkey-adapted *Plasmodium vivax* genomes. *PLoS Negl Trop Dis* **9**: e0003566.
- 674 Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012.  
675 A program for annotating and predicting the effects of single nucleotide polymorphisms,  
676 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*  
677 (*Austin*) **6**: 80-92.
- 678 Claessens A, Affara M, Assefa SA, Kwiatkowski DP, Conway DJ. 2017. Culture adaptation of  
679 malaria parasites selects for convergent loss-of-function mutants. *Sci Rep* **7**: 41303.
- 680 Collaborators GBDDD. 2017. Estimates of global, regional, and national morbidity, mortality, and  
681 aetiologies of diarrhoeal diseases: a systematic analysis for the Global Burden of  
682 Disease Study 2015. *Lancet Infect Dis* **17**: 909-948.
- 683 Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. 2005. Blast2GO: a universal  
684 tool for annotation, visualization and analysis in functional genomics research.  
685 *Bioinformatics (Oxford, England)* **21**: 3674-3676.
- 686 Darriba D, Posada D, Kozlov AM, Stamatakis A, Morel B, Flouri T. 2020. ModelTest-NG: A New  
687 and Scalable Tool for the Selection of DNA and Protein Evolutionary Models. *Mol Biol*  
688 *Evol* **37**: 291-294.
- 689 de Meeus T, Michalakis Y, Renaud F. 1998. Santa rosalia revisited: or why are there so many  
690 kinds of parasites in 'the garden of earthly delights'? *Parasitol Today* **14**: 10-13.
- 691 DeCicco RePass MA, Chen Y, Lin Y, Zhou W, Kaplan DL, Ward HD. 2017. Novel Bioengineered  
692 Three-Dimensional Human Intestinal Model for Long-Term Infection of *Cryptosporidium*  
693 *parvum*. *Infect Immun* **85**.
- 694 Diaz-Viraque F, Pita S, Greif G, de Souza RCM, Iraola G, Robello C. 2019. Nanopore  
695 Sequencing Significantly Improves Genome Assembly of the Protozoan Parasite  
696 *Trypanosoma cruzi*. *Genome Biol Evol* **11**: 1952-1957.
- 697 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome  
698 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* **16**:  
699 157.
- 700 English AC, Salerno WJ, Reid JG. 2014. PBHoney: identifying genomic variants via long-read  
701 discordance and interrupted mapping. *BMC Bioinformatics* **15**: 180.
- 702 Fei J, Wu H, Su J, Jin C, Li N, Guo Y, Feng Y, Xiao L. 2018. Characterization of MEDLE-1, a  
703 protein in early development of *Cryptosporidium parvum*. *Parasit Vectors* **11**: 312.
- 704 Feng Y, Ryan UM, Xiao L. 2018. Genetic Diversity and Population Structure of *Cryptosporidium*.  
705 *Trends in Parasitology* **34**: 997-1011.
- 706 Ghaffari S, Kalantari N, C AH. 2014. A Multi-Locus Study for Detection of *Cryptosporidium*  
707 Species Isolated from Calves Population, Liverpool; UK. *Int J Mol Cell Med* **3**: 35-42.
- 708 Gissot M, Choi SW, Thompson RF, Grealley JM, Kim K. 2008. *Toxoplasma gondii* and  
709 *Cryptosporidium parvum* lack detectable DNA cytosine methylation. *Eukaryot Cell* **7**:  
710 537-540.
- 711 Guo Y, Tang K, Rowe LA, Li N, Roellig DM, Knipe K, Frace M, Yang C, Feng Y, Xiao L. 2015.  
712 Comparative genomic analysis reveals occurrence of genetic recombination in virulent  
713 *Cryptosporidium hominis* subtypes and telomeric gene duplications in *Cryptosporidium*  
714 *parvum*. *BMC Genomics* **16**: 320.
- 715 Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUASt: quality assessment tool for genome  
716 assemblies. *Bioinformatics (Oxford, England)* **29**: 1072-1075.
- 717 Hackl T, Hedrich R, Schultz J, Forster F. 2014. proovread: large-scale high-accuracy PacBio  
718 correction through iterative short read consensus. *Bioinformatics (Oxford, England)* **30**:  
719 3004-3011.

- 720 Heiges M, Wang H, Robinson E, Aurrecochea C, Gao X, Kaluskar N, Rhodes P, Wang S, He  
721 CZ, Su Y et al. 2006. *CryptoDB: a Cryptosporidium* bioinformatics resource update.  
722 *Nucleic Acids Res* **34**: D419-422.
- 723 Heo I, Dutta D, Schaefer DA, Iakobachvili N, Artegiani B, Sachs N, Boonekamp KE, Bowden G,  
724 Hendrickx APA, Willems RJL et al. 2018. Modelling *Cryptosporidium* infection in human  
725 small intestinal and lung organoids. *Nat Microbiol* **3**: 814-823.
- 726 Isaza JP, Galvan AL, Polanco V, Huang B, Matveyev AV, Serrano MG, Manque P, Buck GA,  
727 Alzate JF. 2015. Revisiting the reference genomes of human pathogenic  
728 *Cryptosporidium* species: reannotation of *C. parvum* Iowa and a new *C. hominis*  
729 reference. *Sci Rep* **5**: 16324.
- 730 Jaskiewicz JJ, Sandlin RD, Swei AA, Widmer G, Toner M, Tzipori S. 2018. Cryopreservation of  
731 infectious *Cryptosporidium parvum* oocysts. *Nat Commun* **9**: 2883.
- 732 Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ, de Bakker PI. 2008. SNAP: a  
733 web-based tool for identification and annotation of proxy SNPs using HapMap.  
734 *Bioinformatics (Oxford, England)* **24**: 2938-2939.
- 735 Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A,  
736 Nuka G et al. 2014. InterProScan 5: genome-scale protein function classification.  
737 *Bioinformatics (Oxford, England)* **30**: 1236-1240.
- 738 Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015. The Phyre2 web portal for  
739 protein modeling, prediction and analysis. *Nat Protoc* **10**: 845-858.
- 740 Khalil IA, Troeger C, Rao PC, Blacker BF, Brown A, Brewer TG, Colombara DV, De Hostos EL,  
741 Engmann C, Guerrant RL et al. 2018. Morbidity, mortality, and long-term consequences  
742 associated with diarrhoea from *Cryptosporidium* infection in children younger than 5  
743 years: a meta-analysis study. *Lancet Glob Health* **6**: e758-e768.
- 744 Khan A, Shaik JS, Grigg ME. 2017. Genomics and molecular epidemiology of *Cryptosporidium*  
745 species. *Acta tropica* doi:10.1016/j.actatropica.2017.10.023.
- 746 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory  
747 requirements. *Nat Methods* **12**: 357-360.
- 748 Kirkpatrick M, Barton N. 2006. Chromosome inversions, local adaptation and speciation.  
749 *Genetics* **173**: 419-434.
- 750 Kissinger JC, DeBarry J. 2011. Genome cartography: charting the apicomplexan genome.  
751 *Trends Parasitol* **27**: 345-354.
- 752 Kotloff KL, Nataro JP, Blackwelder WC, Nasrin D, Farag TH, Panchalingam S, Wu Y, Sow SO,  
753 Sur D, Breiman RF et al. 2013. Burden and aetiology of diarrhoeal disease in infants and  
754 young children in developing countries (the Global Enteric Multicenter Study, GEMS): a  
755 prospective, case-control study. *Lancet* **382**: 209-222.
- 756 Krogh A, Larsson B, von Heijne G, Sonnhammer EL. 2001. Predicting transmembrane protein  
757 topology with a hidden Markov model: application to complete genomes. *Journal of*  
758 *molecular biology* **305**: 567-580.
- 759 Lecomte V, Chumpitazi BF, Pasquier B, Ambroise-Thomas P, Santoro F. 1992. Brain-tissue  
760 cysts in rats infected with the RH strain of *Toxoplasma gondii*. *Parasitol Res* **78**: 267-  
761 269.
- 762 Lee E, Helt GA, Reese JT, Munoz-Torres MC, Childers CP, Buels RM, Stein L, Holmes IH, Elsik  
763 CG, Lewis SE. 2013. Web Apollo: a web-based genomic annotation editing platform.  
764 *Genome Biol* **14**: R93.
- 765 Li B, Wu H, Li N, Su J, Jia R, Jiang J, Feng Y, Xiao L. 2017. Preliminary Characterization of  
766 MEDLE-2, a Protein Potentially Involved in the Invasion of *Cryptosporidium parvum*.  
767 *Front Microbiol* **8**: 1647.
- 768 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.  
769 *Bioinformatics* **25**: 1754-1760.

- 770 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,  
771 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and  
772 SAMtools. *Bioinformatics (Oxford, England)* **25**: 2078-2079.
- 773 Li Y, Baptista RP, Sateriale A, Striepen B, Kissinger JC. 2020. Analysis of Long Non-coding RNA  
774 in *Cryptosporidium parvum* Reveals Significant Stage-Specific Antisense Transcription.  
775 *Frontiers in Cellular and Infection Microbiology* doi:10.3389/fcimb.2020.608298.
- 776 Lomsadze A, Ter-Hovhannisyanyan V, Chernoff YO, Borodovsky M. 2005. Gene identification in  
777 novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res* **33**: 6494-6506.
- 778 Mahmoud M, Gobet N, Cruz-Davalos DI, Mounier N, Dessimoz C, Sedlazeck FJ. 2019.  
779 Structural variant calling: the long and the short of it. *Genome Biol* **20**: 246.
- 780 Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, Howard E, Shendure J,  
781 Turner DJ. 2010. Target-enrichment strategies for next-generation sequencing. *Nat*  
782 *Methods* **7**: 111-118.
- 783 Mazurie AJ, Alves JM, Ozaki LS, Zhou S, Schwartz DC, Buck GA. 2013. Comparative genomics  
784 of *Cryptosporidium*. *Int J Genomics* **2013**: 832756.
- 785 McCormick-Graham M, Haynes WJ, Romero DP. 1997. Variable telomeric repeat synthesis in  
786 *Paramecium tetraurelia* is consistent with misincorporation by telomerase. *EMBO J* **16**:  
787 3233-3242.
- 788 McEachern MJ, Haber JE. 2006. Break-induced replication and recombinational telomere  
789 elongation in yeast. *Annu Rev Biochem* **75**: 111-135.
- 790 McEachern MJ, Iyer S. 2001. Short telomeres in yeast are highly recombinogenic. *Mol Cell* **7**:  
791 695-704.
- 792 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytzky A, Garimella K, Altshuler  
793 D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce framework  
794 for analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297-1303.
- 795 Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, Brooks S, Howe E, Porubsky  
796 D, Logsdon GA et al. 2020. Telomere-to-telomere assembly of a complete human X  
797 chromosome. *Nature* **585**: 79-84.
- 798 Morada M, Lee S, Gunther-Cummins L, Weiss LM, Widmer G, Tzipori S, Yarlett N. 2016.  
799 Continuous culture of *Cryptosporidium parvum* using hollow fiber technology. *Int J*  
800 *Parasitol* **46**: 21-29.
- 801 Nader JL, Mathers TC, Ward BJ, Pachebat JA, Swain MT, Robinson G, Chalmers RM, Hunter  
802 PR, Oosterhout C, Tyler KM. 2019. Evolutionary genomics of anthroponosis in  
803 *Cryptosporidium*. *Nature Microbiology*.
- 804 Natarajan S, Nickles K, McEachern MJ. 2006. Screening for telomeric recombination in wild-  
805 type *Kluyveromyces lactis*. *FEMS Yeast Res* **6**: 442-448.
- 806 Nosil P, Feder JL. 2012. Genomic divergence during speciation: causes and consequences.  
807 *Philos Trans R Soc Lond B Biol Sci* **367**: 332-342.
- 808 Otto TD, Sanders M, Berriman M, Newbold C. 2010. Iterative Correction of Reference  
809 Nucleotides (iCORN) using second generation sequencing technology. *Bioinformatics*  
810 **26**: 1704-1707.
- 811 Pawlowic MC, Somepalli M, Sateriale A, Herbert GT, Gibson AR, Cuny GD, Hedstrom L,  
812 Striepen B. 2019. Genetic ablation of purine salvage in *Cryptosporidium parvum* reveals  
813 nucleotide uptake from the host cell. *Proc Natl Acad Sci U S A* **116**: 21160-21165.
- 814 Perteau M, Perteau GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie  
815 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*  
816 **33**: 290-295.
- 817 Piper MB, Bankier AT, Dear PH. 1998. A HAPPY map of *Cryptosporidium parvum*. *Genome Res*  
818 **8**: 1299-1307.
- 819 Ponts N, Fu L, Harris EY, Zhang J, Chung DW, Cervantes MC, Prudhomme J, Atanasova-  
820 Penichon V, Zehraoui E, Bunnik EM et al. 2013. Genome-wide mapping of DNA

- 821           methylation in the human malaria parasite *Plasmodium falciparum*. *Cell Host Microbe*  
822           **14**: 696-706.
- 823 Popendorf K, Tsuyoshi H, Osana Y, Sakakibara Y. 2010. Murasaki: a fast, parallelizable  
824           algorithm to find anchors from multiple genomes. *PLoS One* **5**: e12651.
- 825 Quevillon E, Silventoinen V, Pillai S, Harte N, Mulder N, Apweiler R, Lopez R. 2005.  
826           InterProScan: protein domains identifier. *Nucleic Acids Res* **33**: W116-120.
- 827 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features.  
828           *Bioinformatics (Oxford, England)* **26**: 841-842.
- 829 Ren Q, Chen K, Paulsen IT. 2007. TransportDB: a comprehensive database resource for  
830           cytoplasmic membrane transport systems and outer membrane channels. *Nucleic Acids*  
831           *Res* **35**: D274-279.
- 832 Rieseberg LH. 2001. Chromosomal rearrangements and speciation. *Trends Ecol Evol* **16**: 351-  
833           358.
- 834 Roy A, Kucukural A, Zhang Y. 2010. I-TASSER: a unified platform for automated protein  
835           structure and function prediction. *Nat Protoc* **5**: 725-738.
- 836 Roy A, Yang J, Zhang Y. 2012. COFACTOR: an accurate comparative algorithm for structure-  
837           based protein function annotation. *Nucleic Acids Res* **40**: W471-477.
- 838 Saier MH, Jr., Yen MR, Noto K, Tamang DG, Elkan C. 2009. The Transporter Classification  
839           Database: recent advances. *Nucleic Acids Res* **37**: D274-278.
- 840 Sanderson SJ, Xia D, Prieto H, Yates J, Heiges M, Kissinger JC, Bromley E, Lal K, Sinden RE,  
841           Tomley F et al. 2008. Determining the protein repertoire of *Cryptosporidium parvum*  
842           sporozoites. *Proteomics* **8**: 1398-1414.
- 843 Sateriale A, Slapeta J, Baptista R, Engiles JB, Gullicksrud JA, Herbert GT, Brooks CF, Kugler  
844           EM, Kissinger JC, Hunter CA et al. 2019. A Genetically Tractable, Natural Mouse Model  
845           of Cryptosporidiosis Offers Insights into Host Protective Immunity. *Cell Host Microbe* **26**:  
846           135-146 e135.
- 847 Sateriale A, Striepen B. 2016. Beg, Borrow and Steal: Three Aspects of Horizontal Gene  
848           Transfer in the Protozoan Parasite, *Cryptosporidium parvum*. *PLoS Pathog* **12**:  
849           e1005429.
- 850 Segovia M, Artero JM, Mellado E, Chance ML. 1992. Effects of long-term *in vitro* cultivation on  
851           the virulence of cloned lines of *Leishmania major* promastigotes. *Ann Trop Med Parasitol*  
852           **86**: 347-354.
- 853 Slapeta J. 2013. Cryptosporidiosis and *Cryptosporidium* species in animals and humans: a thirty  
854           colour rainbow? *Int J Parasitol* **43**: 957-970.
- 855 Stanke M, Morgenstern B. 2005. AUGUSTUS: a web server for gene prediction in eukaryotes  
856           that allows user-defined constraints. *Nucleic Acids Res* **33**: W465-467.
- 857 Strong WB, Gut J, Nelson RG. 2000. Cloning and sequence analysis of a highly polymorphic  
858           *Cryptosporidium parvum* gene encoding a 60-kilodalton glycoprotein and  
859           characterization of its 15- and 45-kilodalton zoite surface antigen products. *Infection and*  
860           *Immunity* **68**: 4117-4134.
- 861 Sutherland IA, Shiels BR, Jackson L, Brown DJ, Brown CG, Preston PM. 1996. *Theileria*  
862           *annulata*: altered gene expression and clonal selection during continuous *in vitro* culture.  
863           *Exp Parasitol* **83**: 125-133.
- 864 Swain MT, Tsai IJ, Assefa SA, Newbold C, Berriman M, Otto TD. 2012. A post-assembly  
865           genome-improvement toolkit (PAGIT) to obtain annotated genomes from contigs. *Nat*  
866           *Protoc* **7**: 1260-1284.
- 867 Tandel J, English ED, Sateriale A, Gullicksrud JA, Beiting DP, Sullivan MC, Pinkston B, Striepen  
868           B. 2019. Life cycle progression and sexual development of the apicomplexan parasite  
869           *Cryptosporidium parvum*. *Nat Microbiol* **4**: 2226-2236.

- 870 Troell K, Hallstrom B, Divne AM, Alsmark C, Arrighi R, Huss M, Beser J, Bertilsson S. 2016.  
871 *Cryptosporidium* as a testbed for single cell genome characterization of unicellular  
872 eukaryotes. *BMC Genomics* **17**: 471.
- 873 Vembar SS, Seetin M, Lambert C, Nattestad M, Schatz MC, Baybayan P, Scherf A, Smith ML.  
874 2016. Complete telomere-to-telomere de novo assembly of the *Plasmodium falciparum*  
875 genome through long-read (>11 kb), single molecule, real-time sequencing. *DNA Res*  
876 **23**: 339-351.
- 877 Vinayak S, Pawlowic MC, Sateriale A, Brooks CF, Studstill CJ, Bar-Peled Y, Cipriano MJ,  
878 Striepen B. 2015. Genetic modification of the diarrhoeal pathogen *Cryptosporidium*  
879 *parvum*. *Nature* **523**: 477-480.
- 880 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,  
881 Wortman J, Young SK et al. 2014. Pilon: an integrated tool for comprehensive microbial  
882 variant detection and genome assembly improvement. *PLoS One* **9**: e112963.
- 883 Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B, Guo H et al. 2012.  
884 MCSanX: a toolkit for detection and evolutionary analysis of gene synteny and  
885 collinearity. *Nucleic Acids Res* **40**: e49.
- 886 Wilke G, Funkhouser-Jones LJ, Wang Y, Ravindran S, Wang Q, Beatty WL, Baldrige MT,  
887 VanDussen KL, Shen B, Kuhlenschmidt MS et al. 2019. A Stem-Cell-Derived Platform  
888 Enables Complete *Cryptosporidium* Development In Vitro and Genetic Tractability. *Cell*  
889 *Host Microbe* **26**: 123-134 e128.
- 890 Zahedi A, Monis P, Aucote S, King B, Paparini A, Jian F, Yang R, Oskam C, Ball A, Robertson I  
891 et al. 2016. Zoonotic *Cryptosporidium* Species in Animals Inhabiting Sydney Water  
892 Catchments. *PLoS One* **11**: e0168169.
- 893 Zhang S, Wang Y, Wu H, Li N, Jiang J, Guo Y, Feng Y, Xiao L. 2019. Characterization of a  
894 Species-Specific Insulinase-Like Protease in *Cryptosporidium parvum*. *Front Microbiol*  
895 **10**: 354.  
896