



Long-read assembly and comparative evidence-based reanalysis of *Cryptosporidium* genome sequences reveals expanded transporter repertoire and duplication of entire chromosome ends including subtelomeric regions

Rodrigo P. Baptista, Yiran Li, Adam Sateriale, et al.

Genome Res. published online November 11, 2021
Access the most recent version at doi:[10.1101/gr.275325.121](https://doi.org/10.1101/gr.275325.121)

P<P	Published online November 11, 2021 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Title:** Long-read assembly and comparative evidence-based reanalysis of *Cryptosporidium*
2 genome sequences reveals expanded transporter repertoire and duplication of entire
3 chromosome ends including subtelomeric regions

4
5 **Authors:** Rodrigo P. Baptista^{1,2}, Yiran Li², Adam Sateriale³, Mandy J. Sanders⁴, Karen L.
6 Brooks⁴, Alan Tracey⁴, Brendan R. E. Ansell⁵, Aaron R. Jex⁵, Garrett W. Cooper⁶, Ethan D.
7 Smith⁶, Rui Xiao², Jennifer E. Dumaine³, Peter Georgeson^{6,7,9}, Bernard J. Pope^{6,7,8,10}, Matthew
8 Berriman⁴, Boris Striepen³, James A. Cotton⁴ and Jessica C. Kissinger^{1,2,11}

9
10 **Affiliation:** ¹Center for Tropical and Emerging Global Diseases; ²Institute of Bioinformatics,
11 University of Georgia, Athens, GA, USA; ³Department of Pathology, School of Veterinary
12 Medicine, University of Pennsylvania, Philadelphia, PA, USA; ⁴The Wellcome Sanger Institute,
13 Hinxton, UK; ⁵Faculty of Veterinary and Agricultural Sciences, The University of Melbourne and
14 Population Health and Immunity Division, the Walter and Eliza Hall Institute of Medical
15 Research, Melbourne, Australia; ⁶Department of Clinical Pathology; ⁷Melbourne Bioinformatics;
16 ⁸Department of Surgery (Royal Melbourne Hospital), Melbourne Medical School, Faculty of
17 Medicine, Dentistry and Health Sciences, The University of Melbourne, Australia;
18 ⁹University of Melbourne Centre for Cancer Research, Victorian Comprehensive Cancer
19 Centre, Melbourne, Australia; ¹⁰Department of Medicine, Central Clinical School, Faculty of
20 Medicine Nursing and Health Sciences, Monash University, Australia; ¹¹Department of
21 Genetics, University of Georgia, Athens, GA, USA

22
23 **Key words:** Reassembly, Long reads, genome rearrangement, subtelomeric duplication,
24 subtelomere, reannotation

25

26 **Corresponding Author:** Jessica C. Kissinger jkissing@uga.edu

27

28 **Manuscript Type:** Research

29

30 **Running Title:** New Comparative *Cryptosporidium* Genomic Insights

31

32 **ABSTRACT**

33 Cryptosporidiosis is a leading cause of waterborne diarrheal disease globally and an important
34 contributor to mortality in infants and the immunosuppressed. Despite its importance, the
35 *Cryptosporidium* community has only had access to a good, but incomplete, *Cryptosporidium*
36 *parvum* IOWA reference genome sequence. Incomplete reference sequences hamper
37 annotation, experimental design and interpretation. We have generated a new *C. parvum* IOWA
38 genome assembly supported by PacBio and Oxford Nanopore long-read technologies and a
39 new comparative and consistent genome annotation for three closely related species *C. parvum*,
40 *Cryptosporidium hominis* and *Cryptosporidium tyzzeri*. We made 1,926 *C. parvum* annotation
41 updates based on experimental evidence. They include new transporters, ncRNAs, introns and
42 altered gene structures. The new assembly and annotation revealed a complete *Dnmt2*
43 methylase ortholog. Comparative annotation between *C. parvum*, *C. hominis* and *C. tyzzeri*
44 revealed that most “missing” orthologs are found suggesting that the biological differences
45 between the species must result from gene copy number variation, differences in gene
46 regulation and single nucleotide variants (SNVs). Using the new assembly and annotation as
47 reference, 190 genes are identified as evolving under positive selection, including many not
48 detected previously. The new *C. parvum* IOWA reference genome assembly is larger, gap free
49 and lacks ambiguous bases. This chromosomal assembly recovers all 16 chromosome ends, 13
50 of which are contiguously assembled. The three remaining chromosome ends are provisionally
51 placed. These ends represent duplication of entire chromosome ends including subtelomeric

52 regions revealing a new level of genome plasticity that will both inform and impact future
53 research.

54

55 INTRODUCTION

56 *Cryptosporidium* spp. are parasitic apicomplexans that cause moderate-to-severe diarrhea in
57 humans and animals. Studies have revealed that *Cryptosporidium* is one of the most common
58 causes of waterborne disease in humans and the second leading cause of diarrheal etiology in
59 children < 2 years (Kotloff et al. 2013; GBDDD-Collaborators 2017). In 2016, acute infections
60 caused > 48,000 global deaths and more than 4.2 million disability-adjusted life years lost
61 (Khalil et al. 2018).

62 Currently, 38 species of *Cryptosporidium* are recognized (Slapeta 2013; Feng et al.
63 2018). Most species have preferred hosts, and hosts range from fish to mammals. 15 species
64 have an assembled genome sequence, however, only 8 are annotated (Supplemental Table S1).
65 Most genomic sequence data are from the zoonotic *C. parvum* and anthroponotic *C. hominis*,
66 the species primarily detected in humans (Chalmers et al. 2011; Zahedi et al. 2016; Khan et al.
67 2017). These two species are only 3-5% divergent at the DNA level (Mazurie et al. 2013).
68 *Cryptosporidium* genome sequences are shorter than most other apicomplexans at around 9
69 Mbp distributed over 8 chromosomes and containing < 4,000 protein-encoding genes in the
70 species examined. Most genome reduction consists of gene and intron loss, intron shortening
71 and very short intergenic regions (Abrahamsen et al. 2004; Xu et al. 2004; Kissinger and
72 DeBarry 2011).

73 As the *Cryptosporidium* field is exploding with new-found interest and much needed
74 breakthroughs in genetics and culturing (Vinayak et al. 2015; Morada et al. 2016; DeCicco
75 RePass et al. 2017; Heo et al. 2018; Wilke et al. 2019), the limitations of existing reference
76 genome sequences need to be addressed. The *C. parvum* IOWA II reference genome
77 sequence, (*CplRef*), was assembled with a limited physical map (Abrahamsen et al. 2004) and

78 a few hundred ESTs, for training gene finders. Genomic, transcriptomic and proteomic work has
79 been lacking due to the obligate quasi-intracellular nature of portions of the parasite's life cycle,
80 the historical lack of a continuous *in vitro* tissue culture system, the parasite's small size relative
81 to host cells and difficult animal models. The physical map for the *CplRef* assembly was
82 generated using a genome-wide HAPpily anchored physical mapping technique (Piper et al.
83 1998; Bankier et al. 2003). Despite the cutting-edge approaches, some regions, especially
84 chromosome ends, lacked support or were poorly resolved. Subsequent whole genome
85 sequencing data remain unassembled or in a large number of contigs.

86 In 2015, the *CplRef* was re-annotated using new RNA-seq evidence and a new *C.*
87 *hominis* sequence from a recent human isolate (UdeA01) was generated (Isaza et al. 2015).
88 Many ambiguities in gene models were improved but the new *C. hominis* UdeA01 genome is
89 fragmented. Incomplete, misassembled (i.e. gapped sequence, indels, frameshifts, compressed
90 repetitive regions, artifactual inversions) and independently annotated reference genome
91 sequences as are discussed in (Guo et al. 2015) can mislead analyses of the differences
92 between isolates and species owing to these artifacts rather than the biology. Comparative
93 analyses require additional assays to confirm indels 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 reference genome sequences with new long-read
96 technologies.

97 Long-read sequence technologies (PacBio and Oxford Nanopore) are becoming an
98 essential tool to close full genome sequence assemblies across the tree of life (Vembar et al.
99 2016; Berna et al. 2018; Miga et al. 2020). They can be used to resolve complex regions such
100 as repetitive content, structural variants (SVs) such as inversions, translocations and
101 duplications, or for use as scaffolding evidence for existing fragmented genome assemblies
102 (Mahmoud et al. 2019). They are proving crucial for completing pathogen genome sequences
103 that are often riddled with large virulence-related gene families that may have been improperly

104 assembled in shorter-read assemblies (Xia et al. 2021). Here we provide a new *de novo* hybrid
105 long-read assembly for *C. parvum* strain IOWA (*CplA*), and new consistent, comparative
106 genome annotations for *CplA*, *C. hominis* 30976 (*Ch30976*) and *C. tyzzeri* UGA55 (*CtUGA55*).
107 The new data were used to assess genome-level species differences and assess rapidly
108 evolving genes.

109

110 RESULTS

111

112 **An improved long-read genome assembly for *Cryptosporidium parvum* (IOWA-ATCC)**

113 The *CplRef* genome assembly, generated in 2004, has only 10 physical gaps of unknown size,
114 but it has 18,558 ambiguous bases and is missing 6 telomeres. Alignment of 54,882,187
115 Illumina 100 bp paired-end reads (Supplemental Fig. S1) to this reference sequence, revealed
116 many regions that had become collapsed during assembly (Supplemental Table S2). To
117 resolve these issues, we generated a new PacBio+Illumina+Nanopore hybrid genome assembly
118 (Table 1; Supplemental Fig. S1) for the *C. parvum* strain IOWA (ATCC[®]PRA-67DQ[™]), *CplA*. To
119 minimize strain variation differences, we performed our analysis on the IOWA strain. However,
120 because there is a 14-yr window of propagation between these two isolates, and
121 cryopreservation has only recently been developed (Jaskiewicz et al. 2018), we modified the
122 strain name to IOWA-ATCC (*CplA*).

123 The new, *CplA* genome assembly is compared to the current *CplRef* sequence and two
124 closely related species with different host preferences and pathogenicity, *C. hominis* (*Ch30976*)
125 and *C. tyzzeri* (*CtUGA55*) (Slapeta 2013; Nader et al. 2019; Sateriale et al. 2019) (Table 1).
126 These particular assemblies were selected because they are the best available. The new *CplA*
127 long-read assembly increases the genome size by 19,939 bases (~152 kb when including new
128 proposed subtelomeric regions) and putatively identifies all 16 telomeres. There are no gaps

129 and no ambiguous bases. As expected, the *CplA* genome sequence has diverged slightly but
130 shares 99.93% average pairwise identity with the 2004 assembly in regions that are comparable
131 (Supplemental Table S3). The main *Cryptosporidium* subtyping marker, the 60 kDa surface
132 protein (*gp60* locus subtype IIa) shows 4 amino acid differences (2 in the serine repeat region)
133 between *CplA* and *CplRef* (Supplemental Fig. S2; Supplemental Methods).

134

135 **Structural differences between the *C. parvum* IOWA assemblies are confirmed**

136 The 2004 *CplRef* genome assembly used Sanger reads combined with available HAPPY-map
137 data to scaffold the contigs. We compared the *CplRef* and *CplA* assemblies to identify potential
138 rearrangements. Inversions and relocations were detected in Chromosomes 2, 4 and 5 (Fig. 1A).
139 These inversions may be previous assembly artifacts or represent genuine differences between
140 the isolates. We investigated the synteny between *CplRef*, *CplA*, *Ch30976* and *CtUGA55* and
141 observed that *C. hominis* and *C. tyzzeri* also share the Chr4 and Chr5 inversions. Examination
142 of the inverted region boundaries in *CplRef* revealed regions of ambiguous nucleotide bases or
143 physical gaps (Fig. 1A). To further investigate, PCR primers were designed to test each
144 possible inversion arrangement in genomic DNA from *C. parvum* KSU-1 strain 2006 and Bunch
145 Grass farms IOWA (*CpBGF*) (Fig. 1B, Supplemental Fig S3, Supplemental Table S4). The
146 results support the revised assembly orientation. Long-read Oxford Nanopore, ONT, data also
147 support the *CplA* assembly (Supplemental Fig. S4). Better assemblies for the other species will
148 be needed to determine the true level of synteny across these species.

149

150 **Consistent structural gene annotation resolves inconsistencies and improves functional** 151 **annotation**

152 We consistently annotated and compared *CplRef*, *CplA*, *Ch30976*, and *CtUGA55* which have >
153 95% genome identity to assess differences in gene content. The new annotation for each

154 species was generated with three *de novo* approaches and evidence-based manual annotation.
155 Curation of the annotation was performed pairwise between each assembly to take full
156 advantage of syntenic regions. Data from one species could be used to assess computational
157 predictions in others. Using this approach, fragments of genes that were previously missing in *C.*
158 *hominis* were identified. This approach resulted in 1,926 gene alterations in *CplA* when
159 compared to *CplRef*, resulting in improved functional annotation. These changes increase the
160 overall number of predicted genes, introns (100% supported by RNA-seq data) and exons
161 (Table 2; Supplemental Results). The average mRNA length increased. These structural fixes
162 led to the repair of the N-terminus of the methylase ortholog, *Dnmt2* (Supplemental Fig. S5) as
163 well as 523 other genes and 113 fragmented genes previously annotated as pseudogenes.

164 *Cryptosporidium* has a very compact genome with 76.88% covered by protein coding
165 sequences, CDS. As a result, RNA-seq data, which is the best evidence for annotation, contains
166 reads that overlap adjacent genes creating false fusions of exons belonging to different genes.
167 Available strand-specific RNA-seq was used to characterize some of these regions but
168 expression data were not available for all predicted genes (87% of the annotated genes were
169 covered), thus, genes of unknown function in close proximity on the same strand remain
170 problematic. The expression data also revealed 3 putative alternative spliced genes
171 (CPATCC_0027530; CPATCC0027960; CPATCC_0035590) and 474 potential non-coding
172 RNAs (ncRNAs) predominantly anti-sense lncRNAs with differential expression as reported in
173 (Li et al. 2020).

174

175 **Comparative analysis reveals few gene content differences between closely related**

176 ***Cryptosporidium* species**

177 There is a cluster of species. *C. parvum*, *C. hominis*, *C. tyzzeri*, *C. meleagridis* and *C. ubiquitum*
178 that are highly syntenic relative to species outside of this cluster. The syntenic species are
179 biologically distinct and largely host-adapted with the main zoonotic exceptions being *C. parvum*

180 and *C. ubiquitum*. A synteny analysis of the clustered species and *Cryptosporidium muris* as an
181 outgroup reveals high synteny (99.4% to 87%) within the cluster and only 4% synteny to *C.*
182 *muris* (Supplemental Fig. S6; Supplemental Tables S5 and S14).

183 The consistent annotation of the species closest to *CplA* (GCA_015245375.1), *Ch30976*
184 (GCA_001483515.1) and *CtUGA55* (GCA_007210665.1), permitted the analysis of differences
185 in CDS content and CNVs. Orthology analysis revealed that 94% of the genes were conserved
186 among all species. Of the 4,008 ortholog groups identified, most gene families were maintained
187 with a similar number of paralogs (max = 6) detected in the same ortholog group, but variation
188 was detected among singletons (Fig. 2A; Supplemental Table S6). Some of these gene
189 differences appear to be unique to a particular species (Supplemental Table S7). Of the 224
190 singletons detected, we observed only 0, 1 and 1 potential truly species-specific genes in *CplA*,
191 *Ch30976* and *CtUGA55*, respectively following manual inspection (Fig. 2B; Fig. 2D). Both
192 species-specific genes are uncharacterized proteins. The remaining 253 singletons are detected
193 but incomplete in the fragmented assemblies of *Ch30976* and *CtUGA55*, appearing as split
194 genes, frame-shifts, missed calls near a gap and missing subtelomeric regions or contig break
195 and putative false gene predictions in small contigs (Fig. 2C). The major protein-encoding gene
196 content differences between these species are gene copy number variations and not gene
197 presence or absence.

198 To identify and assess putatively overly collapsed repetitive regions within the genome
199 assemblies analyzed in this study, i.e. repetitive regions represented by only a single repeat in
200 the assembly, we mapped Illumina reads from *CplA*, to the new *CplA*, *CplRef*, *Ch30976* and
201 *CtUGA55* genome assemblies (Supplemental Table S2; Supplemental Fig. S1). Our pipeline
202 detected 12 compressions of at least 2× read depth and > 100 bp in length in the *CplRef*
203 genome assembly compared to 6 in the new *CplA* assembly. The 6 compressed regions drops
204 to 4 if the three putative new subtelomeric regions proposed in this study are included (See
205 below). The *Ch30976* and *CtUGA55* genome assemblies contain > 20 compressions mostly

206 due to the short reads used to generate these assemblies. The *CplA* collapsed regions have 2
207 hits in regions with gene annotations in Chr1 and Chr2. Both genic regions are composed of
208 rRNA genes, some uncharacterized proteins, GMP synthase, aspartate-ammonia ligase,
209 tryptophan synthase beta and MEDLE genes, all associated with complex subtelomeric regions
210 discussed below. The four intergenic compressions all match small simple repeat regions
211 (Supplemental Table S8).

212

213 **Functional annotation identifies new protein features**

214 Several approaches to assess function were applied including InterProScan and I-TASSER
215 among others (see Methods). As a result, 138 new *C. parvum* protein annotations were
216 generated or modified. The percentage of *CplA* genes annotated as uncharacterized proteins
217 was reduced from 40% to 33% in all reannotated sequences (Supplemental Table S9). Many
218 new features including domain and repeat content were added to 738 previously
219 uncharacterized proteins. 729 predicted *CplA* CDSs have signal peptides and 1,990 have GO
220 assignments (Supplemental Results). 1,414 CDSs were further assessed for confidence using I-
221 TASSER protein structure searches and 1,008 predicted structures were assigned as high-
222 confidence by random forest categorization. 143 previously uncharacterized proteins in *CplRef*
223 were assigned high confidence GO terms.

224

225 **New transporter genes were identified**

226 We further characterized transporter genes using three different prediction methods. A total of
227 152 proteins in *CplA* and *Ch30976* were identified as transporters including 128 confident
228 candidates and 24 putative candidates (Supplemental Table S10). This represents an increase
229 of 53 transporters relative to the *CplRef* GO annotation (CryptoDB v36) (Heiges et al. 2006).
230 Most identifiable transporters are related to purine metabolism, peptidoglycan biosynthesis,

231 oxidative phosphorylation and N-Glycan biosynthesis pathways (Fig. 3). Six translocases were
232 also identified.

233

234 **Entire subtelomeric regions are duplicated**

235 As shown in the read depth coverage analysis and in Table 1, Supplemental Table S2 and
236 Supplemental Figure S1, the new *CplA* assembly was able to recover ~2.3kb cumulative length
237 in collapsed regions relative to *CplRef*. One subtelomeric region on Chr1 in *CplA*, previously
238 reported on Chr5 in *CplRef* (but not linked to Chr5 in the HAPPY Map), still shows signs of
239 sequence compression suggesting that most of the genes present in this region have more than
240 one copy (Fig. 4A, Supplemental Fig. S7). This region reveals at least 13 genes which vary in
241 copy number between different *Cryptosporidium* species (Supplemental Fig. S8). The genes
242 contained in this region are 18S rRNA, 5S rRNA and 28S rRNA, uncharacterized proteins, a
243 GMP synthase, an aspartate-ammonia ligase, tryptophan synthase beta and a cluster of several
244 MEDLE genes. Some of these genes, such as the tryptophan synthase beta and the MEDLE's
245 are the focus of considerable research since they may be related to parasite survival and are
246 potentially involved in parasite invasion, respectively (Sateriale and Striepen 2016; Li et al.
247 2017; Fei et al. 2018). The predicted number of copies of rRNAs and MEDLE's are
248 underrepresented as they also have paralogs on Chr2 and Chr5, respectively. The Illumina
249 pileup of ~1,350 reads on Chr2, positions 681,607 to 686,953 (Supplemental Table S2,
250 Supplemental Fig. S1, is the region where the 18S/28S rRNA gene(s) are located on this
251 chromosome. The five 18S genes are identical and 28S rRNAs have three gaps (Supplemental
252 Fig. S9). Thus, alignment competition explains why the read coverage varies relative to the
253 equivalent 18S/28S rRNA Chr1 pileups. Regions with pileups on inner portions of Chr5, 7 and 8
254 are low complexity regions composed by tandem repeats (Supplemental Table S8). In Chr5, we
255 have one uncharacterized protein (CPATCC_0023030), full of tandem repeats and good RNA-

256 seq support for its expression. On Chr7 and 8, these regions are smaller than 100 bp and do not
257 contain any annotated genes.

258 Since there is an apparent compression in a subtelomeric region assembly with no gaps
259 and good PacBio long-read coverage, we hypothesized that these extra copies might derive
260 from additional copies of this region. The *CplA* assembly was only missing three telomeric
261 regions, both ends of Chr7 and one telomere of Chr8. Using existing PacBio long reads we
262 were able to identify a few reads that extended into rRNA regions on the chromosomes missing
263 telomeres. We attempted re-assembly with only PacBio reads and we could not resolve the
264 missing regions. Thus, we generated very deep (2260×) ONT single molecule reads from
265 *CpBGF*, (ATCC DNA was not available, only 143 SNVs are detected between the strains of
266 which 108 are Indels)(Supplemental Table S11). The ONT reads revealed related, yet unique
267 subtelomeric regions linked to the chromosomes missing their telomeres, in addition to Chr1
268 (Fig. 4B). We found good ONT long-read support for these regions (Supplemental Fig. S7).
269 Each distinct subtelomeric region begins with chromosome-specific sequences followed by a
270 conserved ribosomal RNA cluster which is followed by the duplicated subtelomeric region and
271 telomere. There are many ONT and PacBio reads that link the unique chromosomal regions
272 and the beginning of the subtelomeric gene families but only a few span the entire chromosome
273 end. We also note that there is slight variation observed among the reads for each subtelomeric
274 region distal to the rRNA cluster.

275

276 **New positively selected genes are identified in *C. parvum***

277 The new gapless *CplA* genome assembly and annotation presented an opportunity to revisit the
278 prediction of genes evolving under positive selection in this species. We performed a Single
279 Nucleotide Variant, SNV, analysis using 136 different *C. parvum* WGS data sets obtained from
280 GenBank (Supplemental Table S12) using the new *CplA* assembly and annotation. A total of
281 24,407 positions were found to contain at least one high-confidence bi-allelic variant. Multiallelic

282 calls were removed to guard against mixed infections. The biallelic variants reflect 3,892 genes,
283 342 of which show a π_N/π_S ratio of non-synonymous/synonymous rates of > 1.5 (Supplemental
284 Table S13). Of the 342, 17 genes were previously identified and 145 are classified as
285 uncharacterized proteins, 105 of which are annotated as having a signal peptide or being
286 secreted. All previously identified genes evolving under positive selection were detected,
287 including: Insulinase-like protein (CPATCC_0017080), an uncharacterized secreted protein
288 (CPATCC_0010380), *gp60* (CPATCC_0012540) and others (Strong et al. 2000; Sanderson et
289 al. 2008; Nader et al. 2019; Zhang et al. 2019). Of the top 10 genes by π_N/π_S ratio, nine appear
290 to be new to this study. Gene family members such as MEDLEs, FLGN and SKSR were also
291 detected but, new members of each of these families are identified as also evolving under
292 positive selection. Since the putative new subtelomeric repeats (Fig. 4) were not included in
293 these analysis (they were identified in a different strain) evolution of the MEDLE genes may be
294 an over estimation. A family of WYLE proteins (Sanderson et al. 2008) is also identified as being
295 positively selected.

296

297 **DISCUSSION**

298 The first genome sequence assembly of *Cryptosporidium parvum* IOWA II referred to as *CplRef*
299 here, was excellent given the technology at the time. As a result, the community has relied on
300 this genome assembly and annotation to this day to design their experiments. However, gaps
301 and ambiguous bases remain, and there was little available expression and orthology evidence
302 at the time to facilitate the annotation. We used PacBio and Illumina sequencing technologies to
303 generate a new complete genome assembly of *C. parvum* strain IOWA-ATCC. We then applied
304 de novo computational and evidence-based annotation approaches with manual curation of two
305 additional species to generate consistent annotation that can be used to detect differences
306 between species and strains. *CplA* DNA was not available for Nanopore sequencing or PCR
307 validation of the assembly, so *CpBGF* DNA (which differs by < 200 SNVs) was used instead.

308 However, all results are consistent when strains can be compared, e.g. compressions of *CplA*
309 Chr1 detected with *CplA* Illumina reads are the same when *CpBGF* is used lending strength to
310 the broader applicability of the findings.

311 The first, expected, finding was that the *C. parvum* IOWA strain is continuing to evolve
312 (Cama et al. 2006) as it is maintained by passage through cattle in a few different locations for
313 research use. Some natural *Cryptosporidium* isolates have been propagated in unnatural hosts
314 before sequencing. Thus, potential selection during the move to a non-natural host and
315 subsequent drift in propagated and naturally circulating parasites has led to accumulated
316 differences. This phenomenon has been observed in other protozoan parasites (Akiyoshi et al.
317 2002; Cama et al. 2006; Chan et al. 2015; Isaza et al. 2015). Genomic DNA for the 2004 *CplRef*
318 and *CplA* were obtained from the same source, but many years apart. We note small
319 differences in the *gp60* sequence, and an overall genome average difference of 0.07% in
320 identity (Supplemental Table S3).

321 We detect chromosomal inversions in *CplA* relative to *CplRef* which have also been
322 detected by others (Guo et al. 2015; Isaza et al. 2015). Chromosomal inversions are known to
323 affect rates of adaptation, speciation, and the evolution of chromosomes (Rieseberg 2001; Guo
324 et al. 2015), but they can also represent assembly artifacts. PCR spanning the genomic regions
325 flanking each major inversion in each orientation using genomic DNA from an unsequenced
326 isolate from 2006, *C. parvum* KSU-1 and *CpBGF* from 2019 validated the long-read *CplA*
327 assembly. Since the other species still lack physical evidence for their chromosomal structures,
328 further long-read sequencing or chromosome conformation capture sequencing, such as Hi-C,
329 will be needed to detect and validate species-specific genomic structural variations for the other
330 *Cryptosporidium* species.

331 *C. hominis* and *C. tyzzeri*, which are 95–97% identical at the nucleotide level to *C.*
332 *parvum* show incongruences in annotated genes with respect to the new *CplA* genome
333 assembly. The differences result in part from numerous sequence gaps and a lack of

334 experimental evidence (e.g. RNA-seq data) to facilitate annotation. Assembly gaps can lead to
335 frame-shift artifacts, fragments of genes split on different contigs and missing genes. These
336 differences affect similarity-based analyses such as ortholog detection, giving the impression
337 that some of these partially annotated genes are unique to a species when, they are not. These
338 mis-interpretations can sabotage some experimental designs and analysis (Baptista and
339 Kissinger 2019). The re-annotation of the original assembly had 114 pseudogenes, now
340 reduced to just one. This improvement facilitated ortholog and functional identification of the
341 genes involved. Assembly gap regions are usually complex, with repetitive sequence patterns,
342 or hypervariable regions in the population analyzed and some have high polymorphism rates.
343 False assumptions regarding species-specific genes can affect many downstream analyses
344 including the detection of highly polymorphic loci.

345 In this study we were able to improve the structural and functional annotation for three
346 *Cryptosporidium* species using two different approaches: (i) inclusion of seven full-length
347 stranded cDNA libraries derived from three time points (0h, 24h and 48h post infection (Tandel
348 et al. 2019), covering ~ 90% of the protein-encoding genes in *C. parvum*; and (ii) by using
349 synteny information to construct a consistent genome annotation between three different closely
350 related species. This approach facilitated a new comparative analysis of genome content
351 between species. Our analyses reveal that *C. parvum*, *C. hominis* and *C. tyzzeri* show few
352 differences in gene content for the regions that can be compared. Most differences are related
353 to slight structural variation, such as small translocations and inversions, and copy number
354 variation as revealed by read depth coverage analysis.

355 Apicomplexans have streamlined genomes that range from ~8.5 to ~125 Mbp (Woo et
356 al. 2015) relative to the only sequenced free-living ancestor, *Chromera velia* at, 194 Mbp
357 (Kissinger and DeBarry 2011). *Cryptosporidium* species have among the most compact
358 apicomplexan genomes with ~3,900 protein-encoding genes and ~77% of the genome
359 sequence being protein encoding. They also lack a mitochondrial genome and apicoplast

360 organelle (Keeling 2004) a finding that holds with our deep sequencing. Thus, the higher
361 number of transporters found in our re-analysis makes biological sense and adds to growing
362 work in this area, e.g. *Cryptosporidium* may have adapted a novel type of nucleotide transporter
363 for ATP uptake from the host (Striepen et al. 2004; Pawlowic et al. 2019). The new *CpIA*
364 assembly and annotation reveals a complete ortholog of the *Dnmt2* methylase family. The *C.*
365 *parvum Dnmt2* sequence was previously annotated as truncated and lacking a DNMT-specific
366 motif containing a prolyl-cysteinyl dipeptide (Abrahamsen et al. 2004; Ponts et al. 2013; Isaza et
367 al. 2015). DNMT2 proteins share high sequence and structural similarity with DNA
368 methyltransferases, however they appear to function primarily as RNA methyltransferases in
369 plants and animals (Goll et al. 2006). Substrates for DNMT2 in protozoa remain unclear.

370 The lack of three telomeres in the new *CpIA* long-read assembly was an intriguing result
371 that could be explained by the detection of three putative similar but not identical copies of
372 subtelomeric regions containing genes including tryptophan synthase beta, the MEDLE genes
373 and 18S, 5.8S and 28S rRNAs among others. This finding raises the possibility that
374 *Cryptosporidium* has recombination between telomeres by break-induced replication, like some
375 yeasts (McEachern and Iyer 2001; McEachern and Haber 2006), or telomere maintenance by
376 recombination as is observed in human cancers (Natarajan et al. 2006). Some genes in this
377 region may be essential for parasite survival (Sateriale and Striepen 2016). It is possible, but
378 remains to be proven, that extra or altered copies of these genes may confer an advantage to
379 individual parasites or the population as a whole. In fact, we have not shown that all 4
380 subtelomeric regions coexist in the same cell, but 4× coverage of these sequences are present
381 in the population sequenced. We have support from single ONT reads indicating that this region
382 is detected on 4 different chromosome ends. The ONT reads also prove that these structures
383 are varying within the sequenced *CpBGF* population (Fig. 4), raising the possibility of
384 recombination or gene conversion during sexual reproduction. This subtelomeric plasticity in
385 which transfer or duplication of important gene sequences between homologous and

386 nonhomologous chromosome ends, may affect genetic manipulations of the parasite and their
387 resulting phenotype. Currently, cloning does not exist for *Cryptosporidium* and oocysts, which
388 can be sequenced (Troell et al. 2016), must still be considered a population of four haploid
389 meiotic progeny (sporozoites). Single-cell sporozoite sequencing will facilitate recombination
390 and sub-telomeric plasticity studies but currently is still impossible in the absence of genome
391 amplification.

392 *Cryptosporidium* species are usually typed and characterized by a small number of
393 genetic markers: 18S, COWP, HSP70, and *gp60* (Ghaffari et al. 2014). As shown in this study
394 *gp60* which is a gene evolving under positive selection used for *Cryptosporidium* subtyping
395 characterization, had small differences between *CplRef* and *CplA*. Using a single marker to
396 characterize an obligately sexual organism with 8 chromosomes is problematic. In this study, we
397 confirm an existing group of genes evolving under positive selection and identify 325 additional
398 potential candidates distributed across all 8 chromosomes. Some of these genes belong to
399 gene families so to avoid artifacts only uniquely mapped reads were used for the SNV analysis.
400 The genes identified here can be used to help the community develop additional markers for
401 typing parasite isolates. However, the global diversity of *Cryptosporidium* is yet to be
402 characterized. Only 136 isolates from a small geographic region have been sampled here.
403 Newer techniques such as hybrid capture bait set techniques (Mamanova et al. 2010) are a
404 powerful future alternative to characterize and select *Cryptosporidium* population variants and
405 better characterize genetic diversity.

406 The new *C. parvum* long-read assembly combined with a consistent comparative
407 annotation has proven powerful. The species analyzed here have different host preferences and
408 pathogenicity. Comparisons of previous sequences and annotation suggested numerous gene
409 content differences. However, this systematic study reveals that the primary differences
410 between the zoonotic *C. parvum*, the anthroponotic *C. hominis* and the rodent-infecting *C.*
411 *tyzzeri* are SNVs and CNVs rather than differences in unique gene content. Finally, new findings

412 related to within parasite and/or within population subtelomeric amplification and variation
413 events in *C. parvum* reveal a new level of genome plasticity that will complicate some genetic
414 manipulations and may affect the organisms' phenotype.

415

416 **METHODS**

417

418 **Sample DNA sources**

419 *C. parvum* IOWA-ATCC (*CplA*) DNA from oocysts/sporozoitcs was purchased from the ATCC.
420 The source was the University of Arizona, Sterling Parasitology Laboratory. It is a GP60 subtype
421 (IIa) like the current *C. parvum* IOWA II reference (*CplRef*) genome sequence. *C. parvum* IOWA
422 DNA was also prepared from oocysts obtained in 2018 from Bunch Grass Farms, Deary, ID,
423 referred to as *CpBGF* in this study. *C. parvum* KSU-1 genomic DNA was also prepared in 2006
424 from oocysts obtained from Steve Upton. Public sequence data were accessed from NCBI
425 BioProjects PRJNA252787, PRJEB3213, PRJNA388495 and PRJEB10000. Accession
426 numbers for the 136 *C. parvum* sequences used for evolutionary analysis are detailed in
427 Supplemental Table S12.

428

429 ***Cryptosporidium parvum* IOWA-ATCC sequencing and genome assembly**

430 PacBio RSII and Illumina HiSeq 2000 sequencing were performed at the Wellcome Sanger
431 Institute, UK. The *CplA* reads were first assembled using the PacBio open source SMRTlink
432 v6.0 from 9 PacBio SMRT cells, with ~75× mean genome coverage. The resulting assembly
433 was then submitted to the accuracy improver tool Sprai 0.9.9.23 ([https://sprai-](https://sprai-doc.readthedocs.io/en/latest/index.html)
434 [doc.readthedocs.io/en/latest/index.html](https://sprai-doc.readthedocs.io/en/latest/index.html)) and then gaps were filled using PBJelly 15.24.8
435 (English et al. 2014) with PacBio reads and IMAGE 2.4.1 (Swain et al. 2012) with Illumina reads.
436 A manual inspection and improvement using GAP5 (Bonfield and Whitwham 2010) was applied,

437 and the final scaffolded genome assembly was polished with Illumina reads using iCORN2 0.95
438 (Otto et al. 2010) and Pilon 1.22 (Walker et al. 2014).

439 ONT single-molecule long-read sequencing was performed on DNA from *CpBGF*
440 (ATCC®PRA-67DQ™ was out of stock) following the recommended R9.4.1 flow cell protocol.
441 MinION ONT sequencing was performed at the Georgia Genomics Bioinformatics Core at the
442 University of Georgia, USA, using an R.9.4 flow cell and the Ligation Sequencing kit (SQK-
443 LSK109). The ONT long reads generated > 2000× coverage of the *C. parvum* genome. This
444 high coverage complemented the PacBio data to confirm and resolve several complex regions
445 (Supplemental Methods). The final assembly was submitted along with *CplRef*, *Ch30976* and
446 *CtUGA55* to QUAST v.5.02 (Gurevich et al. 2013) to compare and evaluate the quality of *CplA*.
447 All sequencing statistics are in Supplemental Table S12.

448

449 ***Cryptosporidium* genome reannotation**

450 Genome annotation was generated with: ab initio prediction using GeneMark-ES 4.57
451 (Lomsadze et al. 2005); evidence-trained predictions using SNAP/MAKER (Cantarel et al. 2008;
452 Johnson et al. 2008) and AUGUSTUS (Stanke and Morgenstern 2005). For training, we used
453 publicly available data from each respective species: RNA-seq, ESTs, previously predicted
454 proteins and MassSpec proteomics data when available. In parallel we also generated
455 transcriptome assemblies using HISAT2 v.2.1.0 (Kim et al. 2015) and StringTie v.1.3.4 (Pertea
456 et al. 2015), and non-coding RNA predictions were generated for *C. parvum* as described (Li et
457 al. 2020). Manual curation of all genes in the context of existing molecular evidence was
458 performed using WebApollo2 (Lee et al. 2013).

459 We performed comparative genome annotation using the Artemis Comparison tool
460 17.0.1 (Carver et al. 2005). OrthoFinder v.2.3.7 (Emms and Kelly 2015) was used to detect
461 paralogs, orthologs and singletons. All singletons were then manually compared using
462 MCScanX 0.8 (Wang et al. 2012) and JBrowse (Buels et al. 2016) to verify their uniqueness and

463 assess the contribution of sequence gaps or misassembly to the findings. We considered the
464 following error types: Split genes caused by frameshifts or early stop-codons, lack of stranded
465 RNA-seq to confirm the gene model, and the presence of a gapped region in the genome
466 assembly. All genes that did not fall into one of these categories were identified as unique.

467

468 **Functional annotation**

469 Following structural annotation, the predicted protein sequences were used to search the
470 Swiss-Prot, Trembl and the NCBI non-redundant Protein database with BLASTP using an e-
471 value threshold at the superfamily level of 1×10^{-6} . Protein structure similarity was explored
472 using I-TASSER (Roy et al. 2010) as in (Ansell et al. 2019) and Supplemental Methods.
473 Blast2GO (Conesa et al. 2005) version 4.1.9 was used to assign Enzyme Code (E.C) and
474 Gene Ontology (GO) terms. We compared the existing protein product names to the new
475 functional results. Some structural information, such as protein domain and repeat pattern
476 content were added to some uncharacterized proteins and nomenclature errors were corrected
477 according to the NCBI guidelines.

478

479 **Transporter prediction**

480 Predicted proteins were submitted to four different transporter prediction methods: (i) BLASTP
481 against TCDB (Saier et al. 2009) with a threshold e-value of 1×10^{-5} cutoff; (ii) TMHMM (Server v.
482 2.0) (Krogh et al. 2001) and SignalP (Server 4.1) (Bendtsen et al. 2004) to reduce false
483 positives from the TCDB BLASTP results. Transporter candidates with no transmembrane
484 domains or candidates with only one transmembrane prediction while having signal peptides
485 predicted were removed; (iii) TransAAP (Ren et al. 2007), a Transporter Classification tool at
486 TransportDB v2.0 (<http://www.membranetransport.org/transportDB2/index.html>), was used to
487 provide information about potential transporter identity and substrate; and (iv) a structural proof
488 for candidate transporters using Phyre2.0 (Kelley et al. 2015). Final candidate transporters were

489 checked according to above results as well as annotations obtained from InterProScan 5.44
490 (Jones et al. 2014).

491

492 **Comparative analyses**

493 Structural variation sites were calculated using NucDiff v2.0.3 (Khelik et al. 2017) and
494 the major inversions observed between *CplRef* and *CplA* were verified by PCR. Primers to test
495 both ends of each orientation were designed using Primer3 v0.4.0 (Untergasser et al. 2012)
496 (Supplemental Table S4). Primers designed to be specific and conserved among the species
497 were tested using an in silico PCR amplification tool (San Millan et al. 2013). These regions
498 contain repeats so the amplicons range from 2kb to 9kb to avoid them. PrimeSTAR GXL DNA
499 polymerase (TAKARA) was used with Long-PCR conditions: initial 3 min hot start at 98 °C, 35
500 cycles of: 10 sec denaturation at 98 °C; 15 sec primer annealing at 55.4 °C; and 10 min
501 elongation at 68 °C; followed by 10 min elongation at 72 °C. PCR products were separated in a
502 0.8% agarose gels and stained with ethidium bromide.

503 The consistency of annotation and potential gene family CNVs, were determined with
504 OrthoFinder v.2.2.7. CNVs were also determined by aligning Illumina sequence reads from
505 each species studied to the new *CplA* genome sequence to check for variations in read
506 depth coverage. Alignment was performed using BWA-MEM 0.7.17 (Li and Durbin 2009) with
507 default options and the alignment depth per base was calculated using BEDTools
508 genomecov 2.29.2 (Quinlan and Hall 2010) and SAMtools depth 1.6 (Li et al. 2009). Read
509 depth coverage plots were generated using the reshape R package (Wickham 2007; R
510 Development Core Team 2011). To avoid the interference of multiply-mapped regions, only
511 mapped reads were kept for this analysis. For plotting purposes the *Ch30976* genome was
512 scaffolded using the *CplA* chromosomal genomic structure using RagTag v.2.0.1 (Alonge et
513 al. 2019).

514

515 **Resolving the structure of repetitive subtelomeric regions**

516 Following the CNV analysis, the sequence content of the subtelomeric compressed regions and
517 their *CplA* assembly non-compressed chromosomal sequence boundaries containing at least
518 ten genes of Chr1, 7 and 8 were used to build a BLAST database. We then used this database
519 and BLASTN 2.10.0 (Camacho et al. 2009) to detect *CpBGF* ONT reads capable of aligning to
520 both subtelomeric and chromosomal boundary regions. The few reads meeting these criteria
521 were evaluated and visualized by aligning all subtelomeric ONT reads to the unique pre-
522 subtelomeric regions of Chromosomes 1, 7 & 8 using the Geneious mapper 2019.1.3
523 (<https://www.geneious.com>) with medium-sensitivity and minimap2 v.2.22 (Li 2018). Finally, the
524 longest ONT reads were polished with Illumina reads and annotated as previously described for
525 gene content analysis.

526

527 **Variant analysis, selection prediction and populational analysis**

528 Illumina sequence reads from 136 different isolates of *C. parvum* from different geographical
529 locations (Supplemental Table S8) as well as *CpBGF* were aligned against the *CplA* reference
530 genome sequence using BWA-MEM. The BAM files were parsed to select uniquely mapped
531 reads using Picard (Broad_Institute) and then submitted to the GATK 3.8 Haplotypecaller
532 (McKenna et al. 2010). The results were filtered by mapping quality > 40 and depth coverage >
533 10. Because mixed infections exist, we restricted analysis to biallelic sites. The individual VCF
534 files were combined into one GVCF file using the GATK tool GenotypeGVCF. After selecting
535 only SNVs from this data, the combined gvcf file was annotated using SnpEff v.4.3 (Cingolani et
536 al. 2012). The SNV variants from the combined annotated gvcf file had their π_N/π_S ratio (Nei
537 and Gojobori 1986) estimated using SNPGenie 1.0 (Nelson et al. 2015). To avoid noise in the
538 data and identify top candidates, genes with ratios > 1.5 were detected and denoted as evolving
539 under positive selection in the *C. parvum* population analyzed. The higher threshold of 1.5 was

540 chosen based on known genes evolving under positive selection, such as *gp60* (Strong et al.
541 2000) and Insulinase-like (Zhang et al. 2019).

542

543 **DATA ACCESS**

544 The sequence data, and annotation generated in this study have been submitted to the NCBI
545 BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number
546 PRJNA573722 and PRJEB3213. Subtelomeric sequences from Chr7 and 8 have GenBank
547 accessions MZ892386, MZ892387 and MZ892388.

548

549 **ACKNOWLEDGMENTS**

550 We would like to thank Dr, Lihua Xiao for sharing *C. hominis* 30976 and permitting us to update
551 the annotation. This work was supported by Bill and Melinda Gates Foundation grant
552 OPP1151701 to JCK, The Wellcome Trust via its core funding of the Wellcome Sanger Institute
553 (grant WT206194) and NHMRC Investigator Grant (APP1194330) to ARJ.

554

555 **AUTHORS CONTRIBUTIONS**

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

561

562 **COMPETING INTEREST STATEMENT**

563 JCK has a financial interest in PacBio.

564

565 **References**

- 566 Abrahamsen MS, Templeton TJ, Enomoto S, Abrahante JE, Zhu G, Lancto CA, Deng M, Liu C,
567 Widmer G, Tzipori S et al. 2004. Complete genome sequence of the apicomplexan,
568 *Cryptosporidium parvum*. *Science* **304**: 441-445.
- 569 Akiyoshi DE, Feng X, Buckholt MA, Widmer G, Tzipori S. 2002. Genetic analysis of a
570 *Cryptosporidium parvum* human genotype 1 isolate passaged through different host
571 species. *Infect Immun* **70**: 5670-5675.
- 572 Alonge M, Soyk S, Ramakrishnan S, Wang X, Goodwin S, Sedlazeck FJ, Lippman ZB, Schatz
573 MC. 2019. RaGOO: fast and accurate reference-guided scaffolding of draft genomes.
574 *Genome Biol* **20**: 224.
- 575 Ansell BRE, Pope BJ, Georgeson P, Emery-Corbin SJ, Jex AR. 2019. Annotation of the *Giardia*
576 proteome through structure-based homology and machine learning. *Gigascience* **8**.
- 577 Bankier AT, Spriggs HF, Fartmann B, Konfortov BA, Madera M, Vogel C, Teichmann SA, Ivens
578 A, Dear PH. 2003. Integrated mapping, chromosomal sequencing and sequence analysis
579 of *Cryptosporidium parvum*. *Genome Res* **13**: 1787-1799.
- 580 Baptista RP, Kissinger JC. 2019. Is reliance on an inaccurate genome sequence sabotaging your
581 experiments? *PLoS Pathog* **15**: e1007901.
- 582 Bendtsen JD, Nielsen H, von Heijne G, Brunak S. 2004. Improved prediction of signal peptides:
583 SignalP 3.0. *J Mol Biol* **340**: 783-795.
- 584 Berna L, Rodriguez M, Chiribao ML, Parodi-Talice A, Pita S, Rijo G, Alvarez-Valin F, Robello
585 C. 2018. Expanding an expanded genome: long-read sequencing of *Trypanosoma cruzi*.
586 *Microb Genom* **4**.
- 587 Bonfield JK, Whitwham A. 2010. Gap5--editing the billion fragment sequence assembly.
588 *Bioinformatics* **26**: 1699-1703.
- 589 Broad_Institute. Picard Tools.
- 590 Buels R, Yao E, Diesh CM, Hayes RD, Munoz-Torres M, Helt G, Goodstein DM, Elsik CG,
591 Lewis SE, Stein L et al. 2016. JBrowse: a dynamic web platform for genome
592 visualization and analysis. *Genome Biol* **17**: 66.
- 593 Cama VA, Arrowood MJ, Ortega YR, Xiao L. 2006. Molecular Characterization of the
594 *Cryptosporidium parvum* IOWA Isolate Kept in Different Laboratories. *J Eukaryot*
595 *Microbiol* **53 Suppl 1**: S40-42.
- 596 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
597 BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 421.
- 598 Cantarel BL, Korf I, Robb SM, Parra G, Ross E, Moore B, Holt C, Sanchez Alvarado A, Yandell
599 M. 2008. MAKER: an easy-to-use annotation pipeline designed for emerging model
600 organism genomes. *Genome Res* **18**: 188-196.
- 601 Carver TJ, Rutherford KM, Berriman M, Rajandream MA, Barrell BG, Parkhill J. 2005. ACT:
602 the Artemis Comparison Tool. *Bioinformatics* **21**: 3422-3423.
- 603 Chalmers RM, Smith R, Elwin K, Clifton-Hadley FA, Giles M. 2011. Epidemiology of
604 anthroponotic and zoonotic human cryptosporidiosis in England and Wales, 2004-2006.
605 *Epidemiol Infect* **139**: 700-712.
- 606 Chan ER, Barnwell JW, Zimmerman PA, Serre D. 2015. Comparative analysis of field-isolate
607 and monkey-adapted *Plasmodium vivax* genomes. *PLoS Negl Trop Dis* **9**: e0003566.
- 608 Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM. 2012.
609 A program for annotating and predicting the effects of single nucleotide polymorphisms,
610 SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly*
611 (*Austin*) **6**: 80-92.

- 612 Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. 2005. Blast2GO: a universal
613 tool for annotation, visualization and analysis in functional genomics research.
614 *Bioinformatics* **21**: 3674-3676.
- 615 DeCicco RePass MA, Chen Y, Lin Y, Zhou W, Kaplan DL, Ward HD. 2017. Novel
616 Bioengineered Three-Dimensional Human Intestinal Model for Long-Term Infection of
617 *Cryptosporidium parvum*. *Infect Immun* **85**.
- 618 Emms DM, Kelly S. 2015. OrthoFinder: solving fundamental biases in whole genome
619 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* **16**: 157.
- 620 English AC, Salerno WJ, Reid JG. 2014. PBHoney: identifying genomic variants via long-read
621 discordance and interrupted mapping. *BMC Bioinformatics* **15**: 180.
- 622 Fei J, Wu H, Su J, Jin C, Li N, Guo Y, Feng Y, Xiao L. 2018. Characterization of MEDLE-1, a
623 protein in early development of *Cryptosporidium parvum*. *Parasit Vectors* **11**: 312.
- 624 Feng Y, Ryan UM, Xiao L. 2018. Genetic Diversity and Population Structure of *Cryptosporidium*.
625 *Trends in Parasitology* **34**: 997-1011.
- 626 GBDDD-Collaborators. 2017. Estimates of global, regional, and national morbidity, mortality,
627 and aetiologies of diarrhoeal diseases: a systematic analysis for the Global Burden of
628 Disease Study 2015. *Lancet Infect Dis* **17**: 909-948.
- 629 Ghaffari S, Kalantari N, C AH. 2014. A Multi-Locus Study for Detection of *Cryptosporidium*
630 Species Isolated from Calves Population, Liverpool; UK. *Int J Mol Cell Med* **3**: 35-42.
- 631 Goll MG, Kirpekar F, Maggert KA, Yoder JA, Hsieh CL, Zhang X, Golic KG, Jacobsen SE,
632 Bestor TH. 2006. Methylation of tRNA^{Asp} by the DNA methyltransferase homolog
633 Dnmt2. *Science* **311**: 395-398.
- 634 Guo Y, Tang K, Rowe LA, Li N, Roellig DM, Knipe K, Frace M, Yang C, Feng Y, Xiao L. 2015.
635 Comparative genomic analysis reveals occurrence of genetic recombination in virulent
636 *Cryptosporidium hominis* subtypes and telomeric gene duplications in *Cryptosporidium*
637 *parvum*. *BMC Genomics* **16**: 320.
- 638 Gurevich A, Saveliev V, Vyahhi N, Tesler G. 2013. QUASt: quality assessment tool for genome
639 assemblies. *Bioinformatics* **29**: 1072-1075.
- 640 Heiges M, Wang H, Robinson E, Aurrecochea C, Gao X, Kaluskar N, Rhodes P, Wang S, He CZ,
641 Su Y et al. 2006. CryptoDB: a *Cryptosporidium* bioinformatics resource update. *Nucleic*
642 *Acids Res* **34**: D419-422.
- 643 Heo I, Dutta D, Schaefer DA, Iakobachvili N, Artegiani B, Sachs N, Boonekamp KE, Bowden G,
644 Hendrickx APA, Willems RJL et al. 2018. Modelling *Cryptosporidium* infection in
645 human small intestinal and lung organoids. *Nat Microbiol* **3**: 814-823.
- 646 Isaza JP, Galvan AL, Polanco V, Huang B, Matveyev AV, Serrano MG, Manque P, Buck GA,
647 Alzate JF. 2015. Revisiting the reference genomes of human pathogenic *Cryptosporidium*
648 species: reannotation of *C. parvum* Iowa and a new *C. hominis* reference. *Sci Rep* **5**:
649 16324.
- 650 Jaskiewicz JJ, Sandlin RD, Swei AA, Widmer G, Toner M, Tzipori S. 2018. Cryopreservation of
651 infectious *Cryptosporidium parvum* oocysts. *Nat Commun* **9**: 2883.
- 652 Johnson AD, Handsaker RE, Pulit SL, Nizzari MM, O'Donnell CJ, de Bakker PI. 2008. SNAP: a
653 web-based tool for identification and annotation of proxy SNPs using HapMap.
654 *Bioinformatics* **24**: 2938-2939.
- 655 Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A,
656 Nuka G et al. 2014. InterProScan 5: genome-scale protein function classification.
657 *Bioinformatics* **30**: 1236-1240.

- 658 Keeling PJ. 2004. Reduction and compaction in the genome of the apicomplexan parasite
659 *Cryptosporidium parvum*. *Dev Cell* **6**: 614-616.
- 660 Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJ. 2015. The Phyre2 web portal for
661 protein modeling, prediction and analysis. *Nat Protoc* **10**: 845-858.
- 662 Khalil IA, Troeger C, Rao PC, Blacker BF, Brown A, Brewer TG, Colombara DV, De Hostos EL,
663 Engmann C, Guerrant RL et al. 2018. Morbidity, mortality, and long-term consequences
664 associated with diarrhoea from *Cryptosporidium* infection in children younger than 5
665 years: a meta-analysis study. *Lancet Glob Health* **6**: e758-e768.
- 666 Khan A, Shaik JS, Grigg ME. 2017. Genomics and molecular epidemiology of *Cryptosporidium*
667 species. *Acta Trop* doi:10.1016/j.actatropica.2017.10.023.
- 668 Khelik K, Lagesen K, Sandve GK, Rognes T, Nederbragt AJ. 2017. NucDiff: in-depth
669 characterization and annotation of differences between two sets of DNA sequences. *BMC*
670 *Bioinformatics* **18**: 338.
- 671 Kim D, Langmead B, Salzberg SL. 2015. HISAT: a fast spliced aligner with low memory
672 requirements. *Nat Methods* **12**: 357-360.
- 673 Kissinger JC, DeBarry J. 2011. Genome cartography: charting the apicomplexan genome. *Trends*
674 *Parasitol* **27**: 345-354.
- 675 Kotloff KL, Nataro JP, Blackwelder WC, Nasrin D, Farag TH, Panchalingam S, Wu Y, Sow SO,
676 Sur D, Breiman RF et al. 2013. Burden and aetiology of diarrhoeal disease in infants and
677 young children in developing countries (the Global Enteric Multicenter Study, GEMS): a
678 prospective, case-control study. *Lancet* **382**: 209-222.
- 679 Krogh A, Larsson B, von Heijne G, Sonnhammer EL. 2001. Predicting transmembrane protein
680 topology with a hidden Markov model: application to complete genomes. *J Mol Biol* **305**:
681 567-580.
- 682 Lee E, Helt GA, Reese JT, Munoz-Torres MC, Childers CP, Buels RM, Stein L, Holmes IH,
683 Elsik CG, Lewis SE. 2013. Web Apollo: a web-based genomic annotation editing
684 platform. *Genome Biol* **14**: R93.
- 685 Li B, Wu H, Li N, Su J, Jia R, Jiang J, Feng Y, Xiao L. 2017. Preliminary Characterization of
686 MEDLE-2, a Protein Potentially Involved in the Invasion of *Cryptosporidium parvum*.
687 *Front Microbiol* **8**: 1647.
- 688 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics (Oxford,*
689 *England)* **34**: 3094-3100.
- 690 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
691 *Bioinformatics* **25**: 1754-1760.
- 692 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
693 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and
694 SAMtools. *Bioinformatics* **25**: 2078-2079.
- 695 Li Y, Baptista RP, Sateriale A, Striepen B, Kissinger JC. 2020. Analysis of Long Non-Coding
696 RNA in *Cryptosporidium parvum* Reveals Significant Stage-Specific Antisense
697 Transcription. *Front Cell Infect Microbiol* **10**: 608298.
- 698 Lomsadze A, Ter-Hovhannisyanyan V, Chernoff YO, Borodovsky M. 2005. Gene identification in
699 novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Res* **33**: 6494-6506.
- 700 Mahmoud M, Gobet N, Cruz-Davalos DI, Mounier N, Dessimoz C, Sedlazeck FJ. 2019.
701 Structural variant calling: the long and the short of it. *Genome Biol* **20**: 246.

- 702 Mamanova L, Coffey AJ, Scott CE, Kozarewa I, Turner EH, Kumar A, Howard E, Shendure J,
703 Turner DJ. 2010. Target-enrichment strategies for next-generation sequencing. *Nat*
704 *Methods* **7**: 111-118.
- 705 Mazurie AJ, Alves JM, Ozaki LS, Zhou S, Schwartz DC, Buck GA. 2013. Comparative
706 genomics of *Cryptosporidium*. *Int J Genomics* **2013**: 832756.
- 707 McEachern MJ, Haber JE. 2006. Break-induced replication and recombinational telomere
708 elongation in yeast. *Annu Rev Biochem* **75**: 111-135.
- 709 McEachern MJ, Iyer S. 2001. Short telomeres in yeast are highly recombinogenic. *Mol Cell* **7**:
710 695-704.
- 711 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
712 Altshuler D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce
713 framework for analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297-
714 1303.
- 715 Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, Brooks S, Howe E,
716 Porubsky D, Logsdon GA et al. 2020. Telomere-to-telomere assembly of a complete
717 human X chromosome. *Nature* **585**: 79-84.
- 718 Morada M, Lee S, Gunther-Cummins L, Weiss LM, Widmer G, Tzipori S, Yarlett N. 2016.
719 Continuous culture of *Cryptosporidium parvum* using hollow fiber technology. *Int J*
720 *Parasitol* **46**: 21-29.
- 721 Nader JL, Mathers TC, Ward BJ, Pachebat JA, Swain MT, Robinson G, Chalmers RM, Hunter
722 PR, Oosterhout C, Tyler KM. 2019. Evolutionary genomics of anthroponosis in
723 *Cryptosporidium*. *Nature Microbiology*.
- 724 Natarajan S, Nickles K, McEachern MJ. 2006. Screening for telomeric recombination in wild-
725 type *Kluyveromyces lactis*. *FEMS Yeast Res* **6**: 442-448.
- 726 Nei M, Gojobori T. 1986. Simple methods for estimating the numbers of synonymous and
727 nonsynonymous nucleotide substitutions. *Mol Biol Evol* **3**: 418-426.
- 728 Nelson CW, Moncla LH, Hughes AL. 2015. SNPGenie: estimating evolutionary parameters to
729 detect natural selection using pooled next-generation sequencing data. *Bioinformatics*
730 (*Oxford, England*) **31**: 3709-3711.
- 731 Otto TD, Sanders M, Berriman M, Newbold C. 2010. Iterative Correction of Reference
732 Nucleotides (iCORN) using second generation sequencing technology. *Bioinformatics* **26**:
733 1704-1707.
- 734 Pawlowic MC, Somepalli M, Sateriale A, Herbert GT, Gibson AR, Cuny GD, Hedstrom L,
735 Striepen B. 2019. Genetic ablation of purine salvage in *Cryptosporidium parvum* reveals
736 nucleotide uptake from the host cell. *Proc Natl Acad Sci U S A* **116**: 21160-21165.
- 737 Pertea M, Pertea GM, Antonescu CM, Chang TC, Mendell JT, Salzberg SL. 2015. StringTie
738 enables improved reconstruction of a transcriptome from RNA-seq reads. *Nat Biotechnol*
739 **33**: 290-295.
- 740 Piper MB, Bankier AT, Dear PH. 1998. A HAPPY map of *Cryptosporidium parvum*. *Genome Res*
741 **8**: 1299-1307.
- 742 Ponts N, Fu L, Harris EY, Zhang J, Chung DW, Cervantes MC, Prudhomme J, Atanasova-
743 Penichon V, Zehraoui E, Bunnik EM et al. 2013. Genome-wide mapping of DNA
744 methylation in the human malaria parasite *Plasmodium falciparum*. *Cell Host Microbe* **14**:
745 696-706.
- 746 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
747 features. *Bioinformatics* **26**: 841-842.

- 748 R Development Core Team. 2011. R: A language and environment for statistical computing. R
749 Foundation for Statistical Computing, Vienna, Austria.
- 750 Ren Q, Chen K, Paulsen IT. 2007. TransportDB: a comprehensive database resource for
751 cytoplasmic membrane transport systems and outer membrane channels. *Nucleic Acids*
752 *Res* **35**: D274-279.
- 753 Rieseberg LH. 2001. Chromosomal rearrangements and speciation. *Trends Ecol Evol* **16**: 351-
754 358.
- 755 Roy A, Kucukural A, Zhang Y. 2010. I-TASSER: a unified platform for automated protein
756 structure and function prediction. *Nat Protoc* **5**: 725-738.
- 757 Saier MH, Jr., Yen MR, Noto K, Tamang DG, Elkan C. 2009. The Transporter Classification
758 Database: recent advances. *Nucleic Acids Res* **37**: D274-278.
- 759 San Millan RM, Martinez-Ballesteros I, Rementeria A, Garaizar J, Bikandi J. 2013. Online
760 exercise for the design and simulation of PCR and PCR-RFLP experiments. *BMC Res*
761 *Notes* **6**: 513.
- 762 Sanderson SJ, Xia D, Prieto H, Yates J, Heiges M, Kissinger JC, Bromley E, Lal K, Sinden RE,
763 Tomley F et al. 2008. Determining the protein repertoire of *Cryptosporidium parvum*
764 sporozoites. *Proteomics* **8**: 1398-1414.
- 765 Sateriale A, Slapeta J, Baptista R, Engiles JB, Gullicksrud JA, Herbert GT, Brooks CF, Kugler
766 EM, Kissinger JC, Hunter CA et al. 2019. A Genetically Tractable, Natural Mouse Model
767 of Cryptosporidiosis Offers Insights into Host Protective Immunity. *Cell Host Microbe* **26**:
768 135-146 e135.
- 769 Sateriale A, Striepen B. 2016. Beg, Borrow and Steal: Three Aspects of Horizontal Gene
770 Transfer in the Protozoan Parasite, *Cryptosporidium parvum*. *PLoS Pathog* **12**: e1005429.
- 771 Slapeta J. 2013. Cryptosporidiosis and *Cryptosporidium* species in animals and humans: a thirty
772 colour rainbow? *Int J Parasitol* **43**: 957-970.
- 773 Stanke M, Morgenstern B. 2005. AUGUSTUS: a web server for gene prediction in eukaryotes
774 that allows user-defined constraints. *Nucleic Acids Res* **33**: W465-467.
- 775 Striepen B, Pruijssers AJ, Huang J, Li C, Gubbels MJ, Umejiego NN, Hedstrom L, Kissinger JC.
776 2004. Gene transfer in the evolution of parasite nucleotide biosynthesis. *Proc Natl Acad*
777 *Sci U S A* **101**: 3154-3159.
- 778 Strong WB, Gut J, Nelson RG. 2000. Cloning and sequence analysis of a highly polymorphic
779 *Cryptosporidium parvum* gene encoding a 60-kilodalton glycoprotein and
780 characterization of its 15- and 45-kilodalton zoite surface antigen products. *Infection and*
781 *Immunity* **68**: 4117-4134.
- 782 Swain MT, Tsai IJ, Assefa SA, Newbold C, Berriman M, Otto TD. 2012. A post-assembly
783 genome-improvement toolkit (PAGIT) to obtain annotated genomes from contigs. *Nat*
784 *Protoc* **7**: 1260-1284.
- 785 Tandel J, English ED, Sateriale A, Gullicksrud JA, Beiting DP, Sullivan MC, Pinkston B,
786 Striepen B. 2019. Life cycle progression and sexual development of the apicomplexan
787 parasite *Cryptosporidium parvum*. *Nat Microbiol* **4**: 2226-2236.
- 788 Troell K, Hallstrom B, Divne AM, Alsmark C, Arrighi R, Huss M, Beser J, Bertilsson S. 2016.
789 *Cryptosporidium* as a testbed for single cell genome characterization of unicellular
790 eukaryotes. *BMC Genomics* **17**: 471.
- 791 Untergasser A, Cutcutache I, Koressaar T, Ye J, Faircloth BC, Remm M, Rozen SG. 2012.
792 Primer3--new capabilities and interfaces. *Nucleic Acids Res* **40**: e115.

- 793 Vembar SS, Seetin M, Lambert C, Nattestad M, Schatz MC, Baybayan P, Scherf A, Smith ML.
794 2016. Complete telomere-to-telomere de novo assembly of the *Plasmodium falciparum*
795 genome through long-read (>11 kb), single molecule, real-time sequencing. *DNA Res* **23**:
796 339-351.
- 797 Vinayak S, Pawlowic MC, Sateriale A, Brooks CF, Studstill CJ, Bar-Peled Y, Cipriano MJ,
798 Striepen B. 2015. Genetic modification of the diarrhoeal pathogen *Cryptosporidium*
799 *parvum*. *Nature* **523**: 477-480.
- 800 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
801 Wortman J, Young SK et al. 2014. Pilon: an integrated tool for comprehensive microbial
802 variant detection and genome assembly improvement. *PLoS One* **9**: e112963.
- 803 Wang Y, Tang H, Debarry JD, Tan X, Li J, Wang X, Lee TH, Jin H, Marler B, Guo H et al. 2012.
804 MCS-X: a toolkit for detection and evolutionary analysis of gene synteny and
805 collinearity. *Nucleic Acids Res* **40**: e49.
- 806 Wickham H. 2007. Reshaping Data with the reshape Package. *Journal of Statistical Software* **21**.
- 807 Wilke G, Funkhouser-Jones LJ, Wang Y, Ravindran S, Wang Q, Beatty WL, Baldrige MT,
808 VanDussen KL, Shen B, Kuhlenschmidt MS et al. 2019. A Stem-Cell-Derived Platform
809 Enables Complete *Cryptosporidium* Development In Vitro and Genetic Tractability. *Cell*
810 *Host Microbe* **26**: 123-134 e128.
- 811 Woo YH, Ansari H, Otto TD, Klinger CM, Kolisko M, Michalek J, Saxena A, Shanmugam D,
812 Tayyrov A, Veluchamy A et al. 2015. Chromerid genomes reveal the evolutionary path
813 from photosynthetic algae to obligate intracellular parasites. *Elife* **4**: e06974.
- 814 Xia J, Venkat A, Bainbridge RE, Reese ML, Le Roch KG, Ay F, Boyle JP. 2021. Third-
815 generation sequencing revises the molecular karyotype for *Toxoplasma gondii* and
816 identifies emerging copy number variants in sexual recombinants. *Genome Res* **31**: 834-
817 851.
- 818 Xu P, Widmer G, Wang Y, Ozaki LS, Alves JM, Serrano MG, Puiu D, Manque P, Akiyoshi D,
819 Mackey AJ et al. 2004. The genome of *Cryptosporidium hominis*. *Nature* **431**: 1107-1112.
- 820 Zahedi A, Monis P, Aucote S, King B, Papparini A, Jian F, Yang R, Oskam C, Ball A, Robertson I
821 et al. 2016. Zoonotic *Cryptosporidium* Species in Animals Inhabiting Sydney Water
822 Catchments. *PLoS One* **11**: e0168169.
- 823 Zhang S, Wang Y, Wu H, Li N, Jiang J, Guo Y, Feng Y, Xiao L. 2019. Characterization of a
824 Species-Specific Insulinase-Like Protease in *Cryptosporidium parvum*. *Front Microbiol*
825 **10**: 354.
826

827 TABLES

Table 1 Comparative *Cryptosporidium* Genome Assembly Statistics

	CplRef	CplA	Ch30976	CtUGA55
Scaffolds	8	8	53	11
Gaps in assembly	10	0	25	97
Total length bp	9,102,324	9,122,263	9,059,225	9,015,713
Compressed regions*	12	6	21	26
Ambiguous nt	18,558	0	1,699	78,408
Telomere #	10	16**	7	8
N50	1,104,417	1,108,396	470,636	1,108,290
GC%	30.23	30.18	30.13	30.25

828 *Number of compressed regions > 100 nt and > 2× average depth.

829 **BioProject PRJNA573722 and sequence records MZ892386, MZ892387 and MZ892388

Table 2 - Reannotation Summary Statistics.

Strains	<i>C. parvum</i> IOWA II			<i>C. hominis</i>		<i>C. tyzzeri</i>
	IOWA II Before*	IOWA II After**	IOWA- ATCC***	UdeA01***	30976***	UGA55***
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	4,424	3,863	3,996	4,037
Number of CDSs	3,805	3,944	3,897	3,818	3,959	3,986
Average CDS length	1,794	1,765	1,799	1,785	1,755	1,735
Number of exons	4,104	5,043	5,800	4,546	5,045	5,136
Number of introns	238	1,020	1,370	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	76.88	76.1	83.6	79.2

830 *Before refers to the 2007 annotation version available from CryptoDB downloads v.35

831 **After refers to the 2018 annotation version submitted by our group available from CryptoDB
832 v.36

833 ***Version of the annotation available in CryptoDB v.50

834

835

836 **FIGURE LEGENDS**

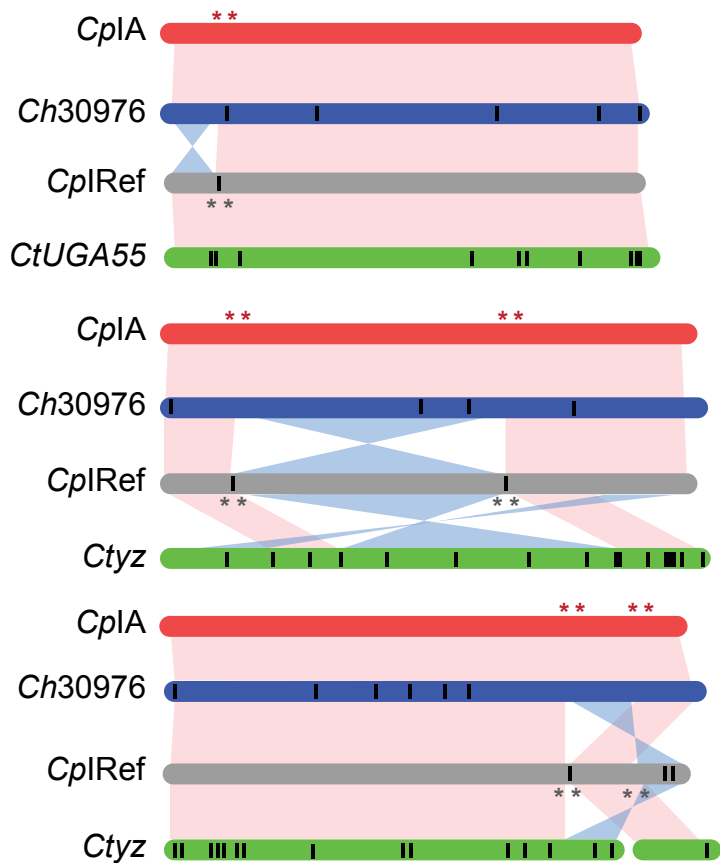
837 **Figure 1.** Syntenic relationships between select *Cryptosporidium* chromosome assemblies. (A)
838 Synteny between Chromosomes 2, 4 and 5. Vertical black lines within a chromosome represent
839 known physical gaps. Synteny between chromosomes is shown in pink and inversions in blue.
840 (B) PCR validation using *C. parvum* KSU-1 DNA, (Supplemental Table S4). Lanes 1, 2 & 3 in all
841 gels are 1kb ladder, positive control *Dnmt2* gene and no template control respectively. The
842 remaining lanes test each orientation of the Left, L, and Right, R, inversion boundaries. Red
843 stars indicate the location of primers designed based on the *CplA* assembly and grey stars
844 indicate the same on the *CplRef* assembly.

845
846 **Figure 2.** Ortholog distribution of protein-encoding genes reveals few differences between the
847 species. (A) Venn diagram of automated protein orthology assignment between *CplA*, *Ch30976*
848 and *CtUGA55*; (B) Venn diagram of the same orthologous genes following manual investigation
849 and removal of putative false positives. *The 139 genes shared between *C. hominis* and *C.*
850 *tyzzeri* in panel A are in complex regions with repeats and gaps and do not have enough
851 evidence to prove their uniqueness at this stage (Supplemental Table S7); (C) Example of a
852 false positive paralog count caused by gene fragments on different scaffolds. (D) Putative
853 unique uncharacterized gene found on Chr8 in *CtUGA55*.

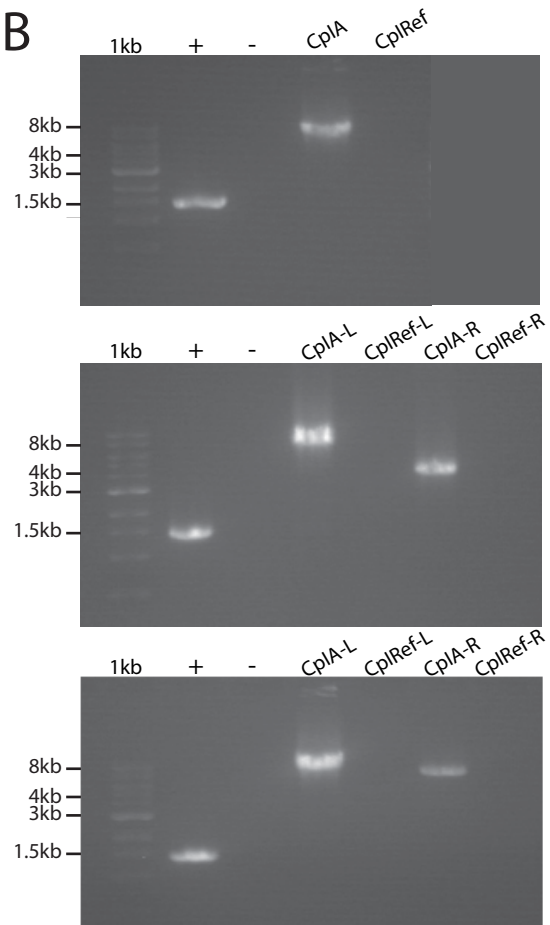
854
855 **Figure 3.** *CplA* assembly and annotation reveal new transporters. Numbers of transporters
856 corresponds to the counts of genes encoding each type of transporter protein. ABC: ATP-
857 binding cassette transporter; MFS: Major facilitator superfamily; DMT: Divalent metal
858 transporter; AAAP: amino acid/auxin permease; MC: mitochondrial carrier; ZIP: Zinc transporter
859 protein; CPA: Cation/Proton Antiporter; SulP: Sulfate Transporter; and PUP: Purine Permeases.
860

861 **Figure 4.** Resolution of repetitive subtelomeric regions found on Chr 1 identifies missing
862 telomeres on Chromosomes 7 and 8. (A) Illumina reads from *CplA* are mapped to the *CplA*
863 Chr1 long-read assembly subtelomeric region to identify read pileups and estimates of
864 sequence copy number by normalizing against the average genomics Illumina read depth.
865 Vertical grey areas indicate regions with annotated genes. Annotated genes are represented
866 below the shaded regions, the 5.8S rRNA is present but not indicated; (B) Subtelomeric variation
867 observed on different *CplA* chromosomes is supported by *CpBGF* ONT long reads. Individual
868 ONT long reads provide evidence of at least four different, yet related, subtelomeric regions that
869 extend into the chromosomes that were missing telomeres in *CplA* (Chr7 and Chr8) in addition
870 to Chr1. The white and black reference bar above each collection of annotated ONT reads
871 identify the resolved subtelomeric regions (white) and linkage to existing assembly (black). The
872 penultimate read on the Chr7 3' end panel indicates a unique region of insertion (nucleotide
873 positions 1191705-1217462). This region contains mostly uncharacterized proteins and two
874 transferases. Each ONT read is annotated as indicated in the Fig. 5A key.

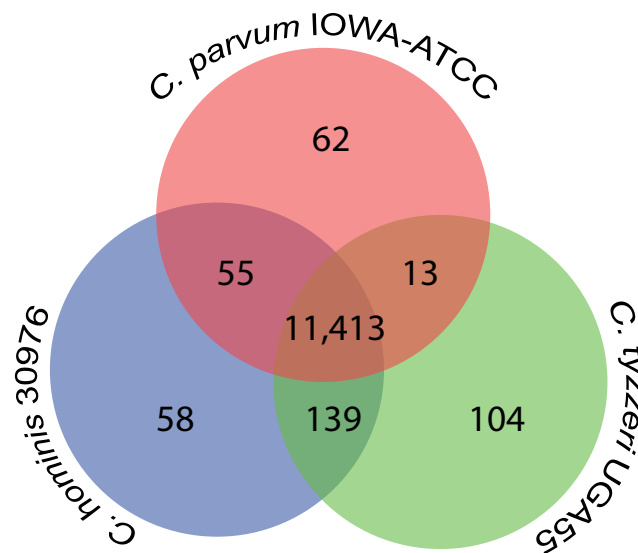
A



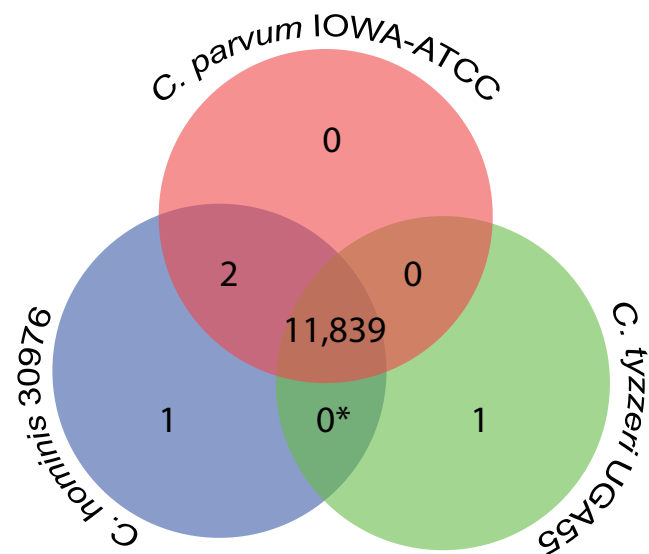
B



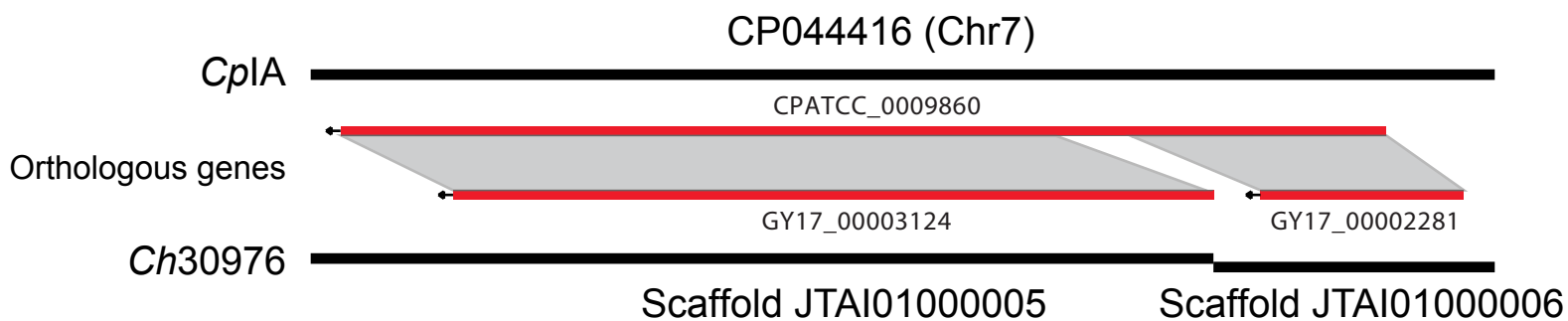
A



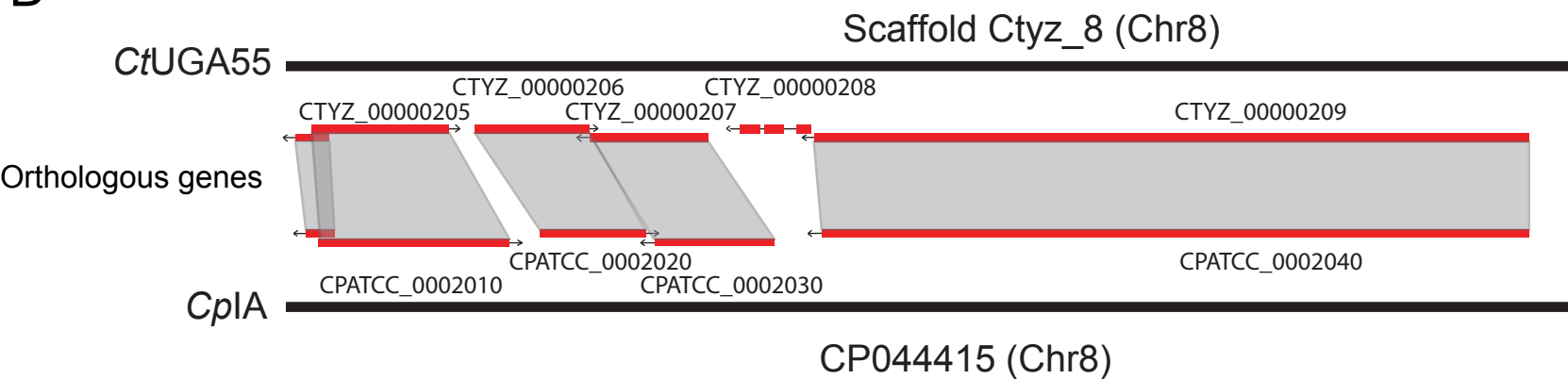
B

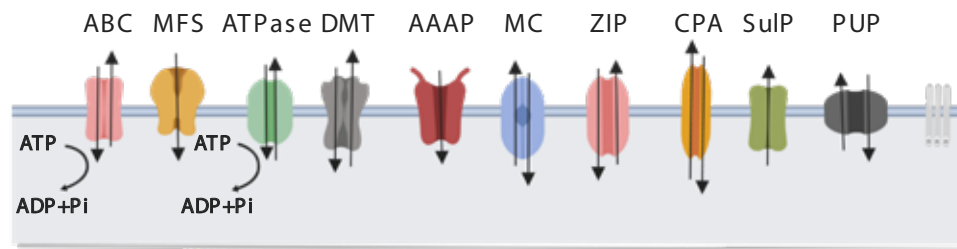


C














D

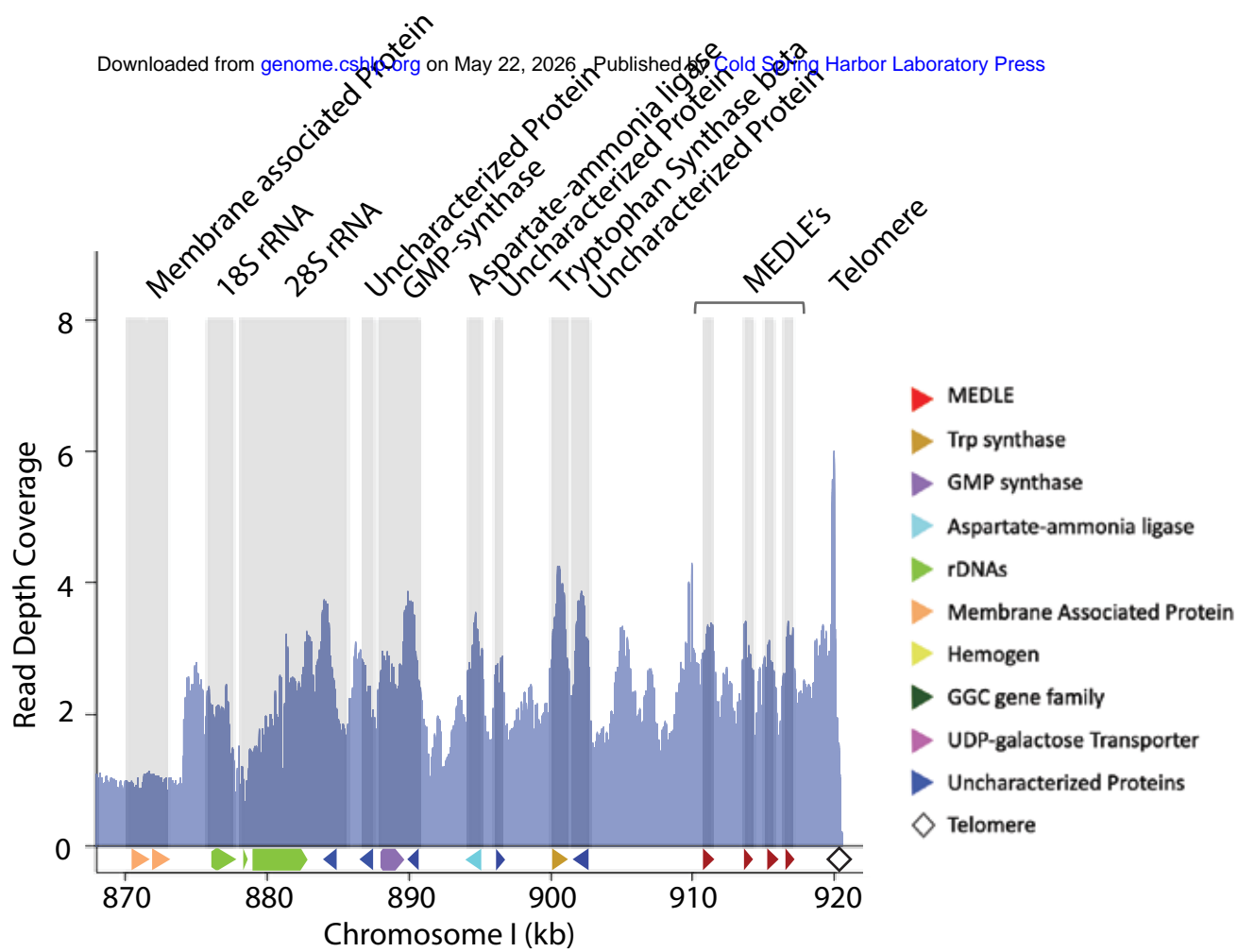




TransportDB	13	7	19	10	7	5	2	1	0	4	16
IOWA-ATCC	19	17	27	13	10	8	3	3	1	4	47

-  ATP-binding cassette transporter
-  Amino acid/auxin permease
-  Sulfate Transporter
-  Major facilitator superfamily
-  Mitochondrial carrier
-  Purine Permeases
-  ATPase transporter
-  Zinc transporter protein
-  Others
-  Divalent metal transporter
-  Cation:Proton Antiporter

A



B

