2

3    **Title**

4    Global population genomics of two subspecies of *Cryptosporidium hominis* during 500 years of

5    evolution

6

7    **Author names**

8    Swapnil Tichkule[1,2,3*], Simone M. Cacciò[4*], Guy Robinson[5,6], Rachel M. Chalmers[5,6], Ivo Mueller[1,3],

9    Samantha J. Emery-Corbin[1], Daniel Eibach[7,8], Kevin M. Tyler[9,10], Cock van Oosterhout[11*] and Aaron

10   R. Jex[1,12*]

11

12   **Affiliation**

13   [1]Population Health and Immunity, Walter and Eliza Hall Institute of Medical Research, Melbourne,

14   VIC, Australia;

15   [2]Faculty of Medicine, Dentistry and Health Sciences, University of Melbourne, Melbourne, VIC,

16   Australia;

17   [3]Department of Medical Biology, University of Melbourne, Melbourne, Australia;

18   [4]Department of Infectious Disease, Istituto Superiore di Sanità, Viale Regina Elena 299, 00161, Rome,

19   Italy;

20   [5]Cryptosporidium Reference Unit, Public Health Wales Microbiology, Singleton Hospital,

21   Swansea, UK;

22   [6]Swansea University Medical School, Swansea, UK;

23   [7]Department of Infectious Disease Epidemiology, Bernhard Nocht Institute for Tropical Medicine

24   Hamburg, Bernhard-Nocht-Strasse 74, 20359 Hamburg, Germany;

25   [8]German Center for Infection Research (DZIF), Hamburg-Lübeck-Borstel-Riems, Germany;

26   [9]Biomedical Research Centre, Norwich Medical School, University of East Anglia, Norwich Research

27   Park, Norwich, UK;

28   [10]Center of Excellence for Bionanoscience Research, King Abdul Aziz University, Jeddah, Saudi

29   Arabia

30   [11]School of Environmental Sciences, University of East Anglia, Norwich Research Park, Norwich, UK;

31   [12]Faculty of Veterinary and Agricultural Sciences, University of Melbourne, Melbourne, VIC, Australia

32

33   **\* Corresponding authors:**

34   Swapnil Tichkule

35    Postal address: Population Health & Immunity, Walter and Eliza Hall Institute of Medical Research,

36    1G, Royal Parade, Parkville VIC 3052

37    Phone: +61 404651636

38    e-mail: tichkule.s@wehi.edu.au

39

40    Simone M. Cacciò

41    Postal address: Department of Infectious Diseases, Istituto Superiore di Sanità, Viale Regina Elena

42    299, Rome 00161 Italy

43    Phone: +39 06 4990 2304

44    e-mail: simone.caccio@iss.it

45

46    Cock van Oosterhout

47    Postal address: School of Environmental Sciences, University of East Anglia, Norwich Research

48    Park, Norwich NR4 7TJ, United Kingdom

49    Phone: +44 (0)1603 592921

50    e-mail: c.van-oosterhout@uea.ac.uk

51

52    Aaron Jex

53    Postal address: Population Health & Immunity Division, Walter & Eliza Hall Institute, 1G Royal

54    Parade, Parkville, VIC 3052

55    Phone: +61 3 9345 2291

56    e-mail: jex.a@wehi.edu.au

57

58    **Abstract**

59    Cryptosporidiosis is a major global health problem and a primary cause of diarrhoea, particularly in

60    young children in low- and middle-income countries (LMICs). The zoonotic *Cryptosporidium parvum*

61    and anthroponotic *C. hominis* cause most human infections. Here, we present a comprehensive whole-

62    genome study of *C. hominis,* comprising 114 isolates from 16 countries within five continents. We

63    detect two lineages with distinct biology and demography, which diverged circa 500 years ago. We

64    consider these lineages two subspecies and propose the names *C. hominis hominis* and *C.*

65    *hominis aquapotentis* (*gp60* subtype IbA10G2). In our study, *C. h. hominis* is almost exclusively

66    represented by isolates from LMICs in Africa and Asia and appears to have undergone recent

67    population contraction. In contrast, *C. h. aquapotentis* was found in high-income countries, mainly in

68    Europe, North America and Oceania, and appears to be expanding. Notably, *C. h. aquapotentis* is

69    associated with high rates of direct human-to-human transmission, which may explain its success in

70    countries with well-developed environmental sanitation infrastructure. Intriguingly, we detected

71 genomic regions of introgression following secondary contact between the subspecies. This resulted

72 in high diversity and divergence in genomic islands of putative virulence genes (GIPVs), including

73 *muc5 (*CHUDEA2_430) and a hypothetical protein (CHUDEA6_5270). This diversity is maintained

74 by balancing selection, suggesting a coevolutionary arms race with the host. Lastly, we find that

75 recent gene flow from *C. h. aquapotentis* to *C. h. hominis*, likely associated with increased human

76 migration, may be driving evolution of more virulent *C. hominis* variants.

77

78 **Introduction**

79 Cryptosporidiosis is a leading cause of diarrhoea in children under five globally (Kotloff, et al. 2013;

80 Khalil, et al. 2018), resulting in an estimated 48,000 deaths annually. Among parasitic diseases, it is

81 second only to malaria in global health burden, with an overall impact of ~12.8 million Disability

82 Adjusted Life-Years (DALYs) (Khalil, et al. 2018). Human cryptosporidiosis is primarily caused by

83 *Cryptosporidium parvum*, a zoonoses common in young ruminants (Ryan, et al. 2016; Santin 2020),

84 and *C. hominis*, which is anthroponotic and the more prevalent species globally (Razakandrainibe, et

85 al. 2018). The disease burden is overwhelmingly skewed to low and middle-income countries

86 (LMICs) (Yang, et al. 2021), particularly in sub-Saharan Africa (Khalil, et al. 2018), where *C.*

87 *hominis* predominates. However, *Cryptosporidium* remains a significant public health problem in

88 wealthy countries through large water or foodborne outbreaks and direct transmission in day-care

89 facilities, hospitals and other institutions (Craun, et al. 1998; Putignani and Menichella 2010).

90 Presently, there are no vaccines or effective drugs to treat the infection. Hence, control depends on

91 prevention of infection, which is driven by a strong understanding of the parasite's epidemiology.

92

93 The global epidemiology of cryptosporidiosis varies by geographic region, socioeconomic

94 status and a range of risk factors (Nichols, et al. 2014; Yang, et al. 2021). Understanding of this

95 epidemiology is underpinned by molecular typing, based mainly on the highly polymorphic *gp60*

96 gene (Feng, et al. 2018). This work has identified numerous genetic variants within each species and

97 indicated a complex population genetic structuring. In *C. parvum*, population structure varies globally

98 from clonal to epidemic to panmictic, likely due to varying ecological factors (Morrison, et al. 2008;

99 Herges, et al. 2012; Wang, et al. 2014). Genetic variants found exclusively in humans point to an

100 anthroponotic *C. parvum* lineage (IIc) (King, et al. 2017; King, et al. 2019), with a recent genomic

101 study recognising two subspecies, the zoonotic *C. parvum parvum* and the anthroponotic *C. parvum*

102 *anthroponosum (*Nader, et al. 2019). Interestingly, this study found that both subspecies occasionally

103 still hybridise and exchange genetic variation. These exchanges overlapped with similar genomic

104 regions undergoing genetic introgression between *C. hominis* and *C. parvum anthroponosum*,

105 indicating candidate sites underpinning adaptation to human-specific infection (Nader, et al. 2019).

106 The signature of introgression can be identified by comparing the DNA sequence (dis)similarity

107 between two or more haplotypes. Briefly, introgressed regions are characterised by high nucleotide

108 similarity, resulting from a relatively recent coalescence of the introgressed sequence. By comparing

109 the sequence variation of three or more haplotypes, the directionality of genetic exchange can be

110 inferred. In addition, the genetic divergence (i.e., number of nucleotide substitutions) can be used to

111 estimate how long ago the sequence was exchanged between the ancestors of the haplotypes. These

112 are the principles of software such as Hybrid-Check (Ward & van Oosterhout 2016) that enable the

113 study of genetic introgression. *Cryptosporidium hominis* appears to have a largely clonal population

114 structure, dominated by specific variants in different regions (Yang, et al. 2021). The IbA10G2 *gp60*

115 subtype, although found in LMICs (Jex and Gasser 2010), is the dominant variant in high-income

116 countries and accounts for up to ~45% of all *gp60*-typed human infections (Jex and Gasser 2010).

117 This subtype is linked to most major waterborne outbreaks in wealthy countries for which genetic

118 typing is available (Zhou, et al. 2003; Chalmers, et al. 2010; Widerström, et al. 2014; Segura, et al.

119 2015; Efstratiou, et al. 2017). The IbA10G2 subtype also appears more readily capable of direct

120 human-to-human transmission (McKerr, et al. 2021) and may cause more severe disease (Cama, et al.

121 2008).

122

123 Studies of *C. hominis* molecular epidemiology pose several essential questions that have

124 major implications for global control of cryptosporidiosis. Specifically: (1) what is the IbA10G2

125 subtype, and is this *gp60*-defined subtype reflective of a phylogenomically divergent *C. hominis*

126 lineage that predominates in wealthy countries; (2) if so, does this lineage undertake reduced levels of

127 genetic recombination with global *C. hominis* populations; (3) and can signatures within the  genome

128 sequences of IbA10G2 typed isolates identify its taxonomic status, reveal the factors underpinning its

129 putatively increased virulence and its dominance in wealthy countries, and identify its influence on

130 global parasite population structure? To address these questions, we performed a global study of 114

131 *C. hominis* genome sequences, comparing the IbA10G2 subtype to published genome sequences

132 representing locally acquired infections from 16 countries across five continents. The insights gained

133 from these analyses are particularly relevant for public health; the IbA10G2 subtype has already been

134 identified as an emerging host-adapted, likely more virulent and transmissible population, making it a

135 threat to human health in high-income countries (Li, et al. 2013; Feng, et al. 2014; Cacciò and

136 Chalmers 2016).

137

138 **Results**

139 **Genomic evidence of population sub-structuring at the continental level**

140 All 114 *C. hominis* isolates included in this study were confirmed as single variant infections

141 (estimated MOI=1 and *Fws* > 0.95) (see Supplementary Table 2). We identified 5,618 biallelic SNPs

142 among these samples and used these to explore the population structure of *C. hominis*. Our analyses

143 identified two major clusters (Fig. 1a), separating European, North American and Oceanian samples

144    from Asian and African samples. We also saw minor clustering separating African from Asian

145    samples. STRUCTURE analysis provided further support for these findings (Fig. 1b). Overall, ~94%

146    of isolates are clustered geographically. Exceptions included a small number of infections (e.g.,

147    UKH30 (from UK) acquired during international travel. STRUCTURE analysis also identified a

148    fourth cluster, including three African and one European isolate, with unique population ancestry

149    (cluster 4 in Fig. 1b).

150

151    **Genomic evidence of population diversification**

152    Maximum Likelihood (ML) analysis (Fig. 1c) identified two major clades, one corresponding to Asia

153    and Africa (clade 1) and the other to Europe, North America and Oceania (clade 2). All isolates

154    within clade 2 had *gp60* subtype IbA10G2, and no IbA10G2 isolates clustered with clade 1. The two

155    clades were estimated to have diverged 488 (84 – 2,199; 95% HPD) years ago. This divergence is

156    supported by a SplitTree network (Fig. 1d) and a DensiTree (Fig. 1e) analyses. Considering the

157    stability of these two genomically distinct lineages and evidence below of their reproductive isolation,

158    we propose their recognition as separate subspecies. We propose clade 1, which comprises infections

159    observed in low- to middle-income countries, mostly from Africa and Asia, be recognised as *C.*

160    *hominis hominis*, referring to the fact that this subspecies represents the majority population. We

161    propose clade 2 (IbA10G2 subtype) includes isolates from high-income countries, namely Europe,

162    North America and Oceania, be named *C. hominis aquapotentis (*strong water), as it predominates in

163    countries with longstanding high sanitation and water quality indices.

164

165    **Demographic histories**

166    To understand the demographic histories and estimate the change in the effective population size (*Ne*)

167    through time, we constructed a Bayesian Skyline Plot (BSP) (Drummond, et al. 2005) for *C. h.*

168    *hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2 IbA10G2). Parasite isolates

169    from low-income countries (*C. h. hominis,* clade 1) experienced a marked population contraction

170    recently from *Ne*=~5000 to *Ne*=~100 (Fig. 2a), which is supported by a higher proportion of positive

171    Tajima's D values (Fig. 2b). In contrast, the isolates from high-income countries (*C. h. aquapotentis*

172    (clade 2), *gp60* subtype IbA10G2) had a stable effective population size (*Ne*=~1000) and a higher

173    proportion of negative Tajima's D values (Fig. 2c). The distribution of Tajima's D in *C. h.*

174    *aquapotentis* is significantly smaller than zero (one-sided t-test: t = -28.883, df = 681, p-value < 2.2e-

175    16), which is consistent with recent population expansion. BSP analysis of *C. h. aquapotentis* (clade

176    2, *gp60* subtype IbA10G2) suggests a stable *Ne*, whereas the overall negative value of Tajima's D

177    analysis is consistent with a recent population size expansion or selective sweep.  BSP analyses are

178     based on coalescent theory, whereas Tajima's D compares the constitution of polymorphisms, i.e.,

179     mean number of pairwise differences and the number of segregating sites. The latter is more sensitive

180     to recent demographic events. Altogether, our analyses thus imply that after a relatively stable *Ne*, the

181     population has started to expand only very recently, or that it has been affected by a recent selective

182     sweep. Simulation-based projections of the evolution of the overall *C. hominis* population finds it is

183     being shaped by 'recent gene flow', likely indicating recent secondary contact (Fig. 2c-d) between the

184     two subspecies ~165 (24 – 647; 5-95% CI) years ago.

185

186          We found evidence of two selective sweeps on chromosome 6 in *C. h. hominis* (clade 1)

187     based on composite likelihood ratio (CLR) statistic using SweeD (Pavlidis, et al. 2013)

188     (Supplementary Fig. 1). However, we do not see any other hallmarks of a selective sweep in these

189     regions based on π or Tajima's D (see Supplementary Fig. 2). Possibly, the selective sweeps were

190     incomplete or occurred in the distant past, eroding their genomic signatures. Nucleotide diversity and

191     Tajima's D are simple statistical approaches to identify selective sweep, but they are also affected by

192     population size changes. In contrast, SweeD analyses site frequency spectra (SFS) and uses complex

193     statistical approaches such as likelihood-based methods to identify selective sweeps in whole genome

194     data. However, the likelihood of selective sweep identified in our analysis is low which might be the

195     reason that Nucleotide diversity and Tajima's D methods did not identify any statistically significant

196     sweeps. Nevertheless, this may have contributed to the recent decline in effective population size of

197     *C. h. hominis* (clade 1) *(*Fig. 2a); genetic variation could have been lost during the selective sweep not

198     only in the affected chromosomal regions, but throughout the genome in this largely clonally

199     reproducing organisms, as the selectively favoured variant replaced other existing variants.

200

201     **Linkage, recombination, introgression and gene flow between subspecies**

202     *Distinct patterns of decay of linkage disequilibria*

203     We performed independent linkage analyses for each subspecies to infer the recombination rate

204     within parasite populations from low- and high-income countries (Fig. 3a). *Cryptosporidium h.*

205     *hominis* (clade 1) had more rapid linkage disequilibrium (LD) decay than *C. h. aquapotentis* (clade 2,

206     *gp60* subtype IbA10G2), consistent with genetic exchanges through gene flow and recombination. In

207     contrast, the strong LD in *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) supports our

208     hypothesis of a recent population expansion (see below).

209

210     *Recombination and regions of secondary contact*

211     Using RDP4 (Martin, et al. 2015), we identified significant recombination events between the two

212     clades in chromosome 1 (event 1), 2 (event 2) and 6 (event 3 and 4) (Fig; 3b and Supplementary

213   Table 3). This indicated secondary contacts between *C. h. hominis* (clade 1) and *C. h. aquapotentis*

214   (clade 2, *gp60* subtype IbA10G2), which resulted in rare genetic exchanges between these otherwise

215   diverged clades. We also undertook additional analyses of the highly admixed European isolate

216   (UK_UKH4 of *gp60* subtype IaA14R3) that showed unique ancestry (Fig 1B) and clustered with low-

217   income countries (*C. h. hominis* (clade 1)) (see Supplementary Text).

218

219   *Signature of introgression and gene flow between the clades*

220   We analysed the recombination events in more detail to better understand the implications of genetic

221   introgression between the two subspecies. Determining the signature of genetic introgression is

222   crucial as these regions can also be responsible for increasing genetic diversity and providing novel

223   substrate for natural selection in host-parasite coevolution (Van Oosterhout 2021). We used

224   HybridCheck (Ward and van Oosterhout 2016) to perform introgression analyses for recombinant

225   events 2 – 4 (event 1 was excluded due to a missing parental sequence). We randomly selected a

226   triplet (recombinant, minor parent and major parental lines) from the RDP output (Supplementary

227   Table 3), which revealed a clear signature of introgression (Fig. 3 c-d), also supported by an ABBA-

228   BABA test (Fig. 3e). Additionally, we calculated the pairwise $r^2$ between SNPs within chromosomes

229   of *C. h. hominis* (clade 1) to assess linkage among SNPs in the introgressed regions (Fig. 3f-g). Large

230   blocks of high LD that encompass the introgressed regions suggested each had been exchanged as a

231   single event between *C. h. hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2)

232   ~165 (24 – 647; 5-95% CI) years ago.

233

234   **Could recent introgression increase virulence?**

235   Our analyses suggest *C. h. hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2)

236   have diverged and largely reproductively isolated for ~500 years. Noting that the latter subspecies has

237   come to dominate infections within high-income countries (Guo, et al. 2015), is more virulent (Cama,

238   et al. 2008), appears better able to transmit through direct human-to-human contact (McKerr, et al.

239   2021) and possibly at a lower infectious dose (Segura, et al. 2015), *C. h. aquapotentis* (clade 2, *gp60*

240   subtype IbA10G2) may owe its success to being better adapted to human infection. If this is the case,

241   it is possible that recent introgression between these two subspecies could select for more virulent and

242   transmissible *C. hominis* subspecies in low-income countries, particularly noting the recent genetic

243   bottlenecking we have observed in *C. h. hominis* (clade 1) within the last decades.

244

245   *Identification of potential virulence genes*

246   To explore this, we first predicted candidate virulence genes in *C. hominis* likely to be involved in the

247   host-parasite interaction and engaged in a coevolutionary arms race with the host. Broadly, such genes

248    are under continuous adaptation while interacting with hosts (Van Oosterhout 2021). During

249    coevolution, evolutionary forces act on virulence genes to create genetic variation through mutation,

250    recombination and gene flow and this variation is moulded by natural selection (and genetic drift). To

251    identify putative virulence genes in *C. hominis*, we selected the top 5% most highly polymorphic

252    genes based on nucleotide diversity ($\pi$). These were filtered for the top 25% of genes under balancing

253    selection, based ranked Tajima's D. Finally, these genes were filtered by selecting the top 50% of

254    genes with the highest proportion of non-synonymous mutations, based on ranked Ka/Ks ratio. Using

255    this approach, we identified 24 highly polymorphic genes (Supplementary Table 5) that are rapidly

256    mutating at the protein level and under selective pressure. These genes were significantly more

257    polymorphic (two-sided t-test: t = -2.9062, df = 23, p-value = 0.007957), under stronger balancing

258    selection (two-sided t-test: t = -10.393, df = 23, p-value = 2.533e-10) and positive selection (two-

259    sided t-test: t = -2.4736, df = 23, p-value = 0.021) compared to all remaining polymorphic genes.

260    Moreover, these genes are enriched in recombination regions ($\chi2$= 225.04, df = 1, p-value = 0.00049),

261    showing that gene flow within *C. hominis* is inclined towards elevated levels of genetic variation.

262    Overall, these genes are enriched for extracellular ($\chi2$= 13.608, df = 1, p-value = 0.00349) and signal

263    peptide proteins ($\chi2$= 6.69, df = 1, p-value = 0.015), which is consistent with their potential relevance

264    for host-parasite interaction. Together, these results provided the evidence of genomic islands of

265    putative virulence genes (GIPVs: Fig. 4) that are likely to be in a coevolutionary arms race with the

266    host, undergoing frequent recombination and accumulating beneficial polymorphisms that are under

267    selection, and that this increased population diversity.

268

269    *Subspecific divergence of putative virulence genes*

270    We investigated whether any virulence genes predicted above were highly diverged between *C. h.*

271    *hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) and under diversifying

272    selection between the two subspecies.  To do this, we calculated absolute divergence (Dxy),

273    representing the average proportion of differences between all pairs of sequences between *C. h.*

274    *hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2), revealing differential gene

275    flow during their reproductive isolation (Fig. 4a). We then calculated the correlation between

276    diversity ($\pi$) and divergence (Dxy) for each gene (Fig. 4b). This approach identified four clear outlier

277    genes, which were the most divergent, diverse, and rapidly mutating at the protein level. These genes

278    encoded three mucins (CHUDEA2_430, CHUDEA2_440 and CHUDEA2_450) arrayed in a cluster

279    on chromosome 2, and a hypothetical protein (CHUDEA6_5270) found on chromosome 6 (Fig. 4b,

280    Supplementary Table 5). The *gp60* gene ranked 6[th] in the list of virulence genes but does not sit as a

281    clear outlier from the other genes identified in our assessment. Two of these genes, CHUDEA2_430

282    (LRT, p-value = 0.005) and CHUDEA6_5270 (LRT, p-value = 0.017), are under statistically

283   significant diversifying selection between *C. h. hominis* (clade 1) and *C. h. aquapotentis* (clade 2,

284   *gp60* subtype IbA10G2).

285         Lastly, we looked at each codon position within each of the four outlier genes to detect codon

286   specific diversifying selection, which might be overlooked in the overall gene. We detected episodic

287   diversifying selection (LRT p-value < 0.05) at codon positions 92, 111 and 138 of CHUDEA2_430

288   and 97 and 105 of CHUDEA6_5270 using the Mixed Effects Model of Evolution (MEME) method

289   implemented in Datamonkey (Pond and Frost 2005). These sites represent putative codons harbouring

290   genetic polymorphisms that experience periods of strong diversifying selection. This strongly

291   suggests CHUDEA2_430 and CHUDEA6_5270 are likely candidate virulence factors participating in

292   a coevolutionary arms race and contributing to the divergence of *C. h. hominis* (clade 1) and *C. h.*

293   *aquapotentis* (clade 2, *gp60* subtype IbA10G2).

294

295   *Recent introgression of putative virulence genes*

296   Finally, we asked whether the recent recombination events in chromosomes 2 and 6 between *C. h.*

297   *hominis* (clade 1) and *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) might include genes

298   associated with increased virulence, indicating increasing virulence of *C. hominis* globally within the

299   past decades, possibly linked to increased human migration. These introgressed regions had

300   significantly elevated levels of nucleotide diversity (two-sided t-test: t = -3.0026, df = 27, p-value =

301   0.0057), divergence (two-sided t-test = -3.0043, df = 27, p-value = 0.0057) and balancing selection

302   (two-sided t-test: t = -3.0125, df = 27, p-value = 0.0055). Interestingly, we observed a pattern of high

303   diversity and divergence (Fig. 4a), particularly for genes within the recombinant blocks (Fig. 3b-e).

304   The 38 (= top 1%) most divergent genes between *C. h. hominis* (clade 1) and *C. h. aquapotentis*

305   (clade 2, *gp60* subtype IbA10G2) were enriched in recombinant regions ($\chi$2=  287.08, df = 1, p-value

306   = 0.00049), with 37% present in these regions (Supplementary Table 6), compared to 0.9% of the

307   other 1,645 divergent genes. These results suggest the genes undergoing frequent recombination

308   accumulated beneficial polymorphisms that were maintained by balancing selection and that this

309   increased population diversity.

310         Notably, the introgressed regions between *C. h. hominis* (clade 1) and *C. h. aquapotentis*

311   (clade 2, *gp60* subtype IbA10G2) included CHUDEA2_430 (*muc5*) for event 2 (chromosome 2),

312   CHUDEA6_1080 (*gp60*) for event 3 (chromosome 6) and CHUDEA6_5270 (a hypothetical protein)

313   for event 4 (chromosome 6). The introgression at CHUDEA6_5270 is particularly intriguing as it

314   sheds further light on the evolution of *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) and *C. h.*

315   *hominis* (clade 1), revealing how recent recombination events could be driving the virulence evolution

316   of *C. h. hominis* (clade 1). We identified two major CHUDEA6_5270 (hypothetical gene) haplotypes;

317   Hap1 representing most *C. h. hominis* (clade 1) isolates, and Hap2 represents all *C. h. aquapotentis*

318  (clade 2, *gp60* subtype IbA10G2) plus a small subset of *C. h. hominis* (clade 1) (Fig. 5a), as well as

319  many haplotypes associated only with *C. h. hominis* (clade 1). No mutations were observed in Hap2,

320  which is consistent with *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) having evolved recently,

321  being an estimated 392 years (29 – 1699 years; 5-95% CI) old, assuming a mutation rate of $u=10^{-8}$

322  and 48h cell division time. Hap1 and Hap2 are diverged by 22 SNPs (Fig. 5b), which is unlikely to

323  represent standing genetic variation, given this is 3-fold higher than the mean deviation among all

324  other CHUDEA6_5270 haplotypes identified here. Instead, we propose CHUDEA6_5270-Hap2

325  might be an introgressed variant from a highly diverged, unsampled (sub)species that diverged from

326  *C. hominis* around 2.36 million generations ago (assuming a mutation rate of $u=10^{-8}$), which is equal

327  to 12,742 (8,901 – 17,497 years; 5-95% CI) years (assuming 48h/replication). The emergence of

328  CHUDEA6_5270-Hap2 in some *C. h. hominis* (clade 1) isolates (i.e., the red section of Hap2 in Fig.

329  6a) strongly implies that *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) has introgressed (circa

330  400 years ago) into *C. h. hominis* (clade 1) leading to evolution under balancing selection (Fig. 5d).

331

332      We modelled the 3D protein structure of CHUDEA6_5270-Hap1 and Hap2, and compared

333  these to similarly predicted 3D protein structures for the orthologous genes from *C. parvum parvum*

334  and *C. parvum anthropanosum (*Fig. 5c). This modelling identified several conserved alpha-helices

335  that appear to form a coiled-coil domain. CHUDEA6_5270 from *C. h. aquapotentis* (clade 2, *gp60*

336  subtype IbA10G2) encodes mutations near the C-terminal end, resulting in a notable kink that

337  deviates from the other structures. This structural variation overlaps with an increase in Tajima's D

338  values toward the C-terminus of the protein, suggesting the region is under balancing selection (Fig.

339  5d).

340

341      A similar pattern is observed for CHUDEA2_430 (see Supplementary Fig. 6). However,

342  given that both the subspecies are interspersed in the haplotype network, indicating against a specific

343  directionality of introgression. This may indicate that CHUDEA_430 diverged with the divergence of

344  the subspecies and continued to diversify. We were not able to generate a robust 3D structural model

345  for CHUDEA2_430 or its *C. parvum* orthologs. We noted the gene encodes a large intrinsically

346  disordered region (Supplementary Fig. 7) which is a consistent feature of mucin proteins (Carmicheal,

347  et al. 2020), whose structural confirmation is influenced by post-translation glycosylation (Perez-Vilar

348  and Hill 1999).  Noting this, we did identify eight novel glycosylation sites in *C. h. aquapotentis*

349  (clade 2, *gp60* subtype IbA10G2) CHUDEA2_430 haplotypes not found in *C. h. hominis* (clade 1)

350  *(*Supplementary Fig 5). Whether these impact interaction with host proteins is not known, but this

351  would be consistent with glycosylated proteins in other pathogens (Lin, et al. 2020).

352

353  **Discussion**

354    In this study, we examined the evolutionary genomics of a major human parasite, *C. hominis*,

355    studying whole genome sequence data from 114 isolates from 16 countries across five continents. We

356    posed three questions, and we are able to answers these as follows: (1) the *gp60*-defined subtype

357    IbA10G2 is reflective of a phylogenomically divergent *C. hominis* lineage that predominates in

358    wealthy countries; (2) this lineage has experienced significantly reduced levels of genetic

359    recombination with other global *C. hominis* populations, and we could identify only four genetic

360    exchanges between these otherwise diverged clades; (3) the distinct evolutionary trajectory of the

361    IbA10G2 subtype, characterized by rapid population expansion, warrants a distinct taxonomic status

362    as subspecies. We propose the name *C. hominis aquapotentis* (clade 2, *gp60* subtype IbA10G2) to

363    reflect its adaptation to "strong water", i.e., high sanitation and water quality indices.

364    These two subspecies are estimated to have been reproductively isolated for approximately

365    488 (84 – 2,199) years, except for the more recent genetic exchange at four genetic loci. The

366    reproductive isolation coincides with improvements to sanitation in Europe. It is possible that *C. h.*

367    *aquapotentis* (clade 2, *gp60* subtype IbA10G2) evolved specialisations making it better suited to

368    human infection allowing it to be more successful through direct transmission supported by a lower

369    infectious dose. Such adaptations would have allowed *C. h. aquapotentis* (clade 2, *gp60* subtype

370    IbA10G2) to become dominant in higher-income countries where sanitation has reduced the level of

371    environmental transmission. Indeed, epidemiological investigations in the UK identified direct

372    person-to-person transmission as a key pathway for *C. h. aquapotentis* (clade 2, *gp60* subtype

373    IbA10G2) (McKerr, et al. 2021). We hypothesise this may have resulted in its rapid population

374    expansion over the last ~500 years and (partial) reproductive isolation. In contrast, *C. h. hominis*

375    (clade 1) experienced recent reduction in effective population size (*Ne*) within the past few decades.

376    We detected two signatures of recent selective sweep in this subspecies, and we propose that this may

377    have eroded some of the genetic variation, resulting in the marked drop in *Ne* = 5000 to *Ne* = 100 in

378    the past. This would have resulted in significant genetic drift and random allele frequency changes,

379    which could have increased the divergence between *C. h. hominis* (clade 1) and *C. hominis*

380    *aquapotentis* (clade 2, *gp60* subtype IbA10G2) further.

381    Despite increased migration and international travelling in the past decades, our study

382    suggests that there has been relatively little movement between continents for this parasite. This is in

383    contrast to reports for *C. parvum*. Corsi, *et. al.* recently found a higher proportion of admixture and

384    gene flow between *C. parvum* populations and no evidence of population structuring by geographic

385    region (Corsi, et al. 2021). Despite the strong population sub-structuring in *C. hominis*, we found

386    evidence of potential recombination and gene flow between the geographic populations and

387    subspecies. We further investigated and identified the introgressed regions where we detected

388    significant gene flow between the low- and high- income countries. Simulation-based analyses

389    indicated this was most likely explained by 'recent geneflow' (circa 165 years ago). This would

390  appear to be a secondary contact between the two subspecies after recent globalisation, illustrating

391  higher migration rate from high-income to low-income countries, which facilitated gene flow,

392  recombination, population admixture and selective sweep.

393      Genetic exchanges between *C. h. hominis* (clade 1) and *C. hominis aquapotentis* (clade 2,

394  *gp60* subtype IbA10G2) are rare compared to those within *C. parvum parvum* (Corsi, et al. 2021), and

395  their frequency might be more comparable to the rate of sequence exchange between *C. p. parvum*

396  and *C. p. anthroponosum* (Nader, et al. 2019). However, whole genome analyses of more *C. hominis*

397  isolates may detect other recombination events in addition to the four events detected in our study. In

398  addition, further studies may be able to discover the unknown parental sequences associated with

399  recombination event 1 in our study.  Without this parental sequence, we were unable to reconstruct

400  the evolution of this introgressed sequence on chromosome 1. Yet, this event may be a key player in

401  the evolution of the lineage. We encourage future whole genome studies on *C. hominis*, believing this

402  may shed further light on the incipient speciation of *C. hominis aquapotentis* (clade 2, *gp60* subtype

403  IbA10G2).

404

405      Although our dataset comprises samples across five continents, we only studied *C. hominis* in

406  16 countries in total, which means that we could have missed local gene flow and patterns of

407  population sub-structuring within continents. Although the marked biological differences between *C.*

408  *parvum* and *C. hominis* have been well established (Abrahamsen, et al. 2004), recent population

409  genomic research is demonstrating that also within these species, the population genetics and

410  evolutionary genetics of their subspecies are remarkably distinct. Large datasets and comparative

411  population genomic and phylogenomic analyses (across *Cryptosporidium* species) are warranted to

412  examine the evolutionary genomics of these parasites in more detail.

413

414      Finally, we have discovered genomic islands of putative virulence genes (GIPVs)

415  contributing to population diversification between *C. h. aquapotentis* (clade 2, *gp60* subtype

416  IbA10G2) and *C. hominis hominis* (clade 1). These islands have experienced relatively elevated

417  recombination rate which has enriched nucleotide variation under balancing selection and the

418  acquisition of non-synonymous SNPs, consistent with virulence factors driving host-parasite

419  interactions. Intriguingly, the most significant signals within these analyses are driven by *gp60*, a

420  hypothetical protein (CHUDEA6_5270) and a cluster of mucin-like genes (CHUDEA2_430,

421  CHUDEA2_440 and CHUDEA2_450) found on chromosome 2. These genes are consistently

422  identified as being under selection in the evolution of the *C. hominis* subspecies here and in similar

423  observations made of *C. hominis* in Africa (Tichkule, et al. 2021). Their orthologs are associated with

424  recombination between human-specific *C. parvum anthroponosum* relative to the zoonotic *C. parvum*

425  *parvum* and appear to have driven convergence of the former with *C. hominis* (Nader, et al. 2019).

426    CHUDEA2_430 (*muc5*) and hypothetical protein CHUDEA6_5270 are the most notable,

427    displaying significant diversifying selection between the two subspecies. Broadly, mucins mediate

428    cell-cell interactions (O'Connor, et al. 2009), and modulate infectivity of *Cryptosporidium* sporozoites

429    and merozoites and oocyst production (Cevallos, et al. 2000; O'Connor, et al. 2009). In *C. parvum*,

430    MUC5 is involved in host-cell invasion and an important determinant of host adaptation (O'Connor, et

431    al. 2009) and highly expressed in the first 2 hours of infection in vitro (Lippuner, et al. 2018). MUC5

432    may also play a role in tethering the sporozoite to the oocyst wall (Chatterjee, et al. 2010). Our

433    analyses suggest *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) and *C. h. hominis* (clade 1)

434    *muc5* haplotypes diverged before or with the subspecies and subsequently diversified, which is

435    consistent with prior observations implicating CHUDEA2_430 in the emergence of *C. h. aquapotentis*

436    (clade 2, *gp60* subtype IbA10G2) (Bouzid, et al. 2013; Feng, et al. 2018). This appears to have

437    resulted in the acquisition of novel glycosylation sites within *C. h. aquapotentis* (clade 2, *gp60*

438    subtype IbA10G2) *muc5* haplotypes. We cannot determine the functional consequence of these sites

439    but note that glycosylation sites often mediate the specificity of mucin interactions with host proteins

440    in a variety of pathogens (Lin, et al. 2020). In contrast, CHUDEA6_5270 displays a clear signal for

441    the recent introgression of a novel *C. h. aquapotentis* (clade 2, *gp60* subtype IbA10G2) haplotype into

442    *C. h. hominis* (clade 1) after the divergence of these subspecies. This haplotype has notable,

443    structurally relevant, mutations. Identifying the function of this gene, its potential role in infection and

444    the relevance of the structural variation we have inferred here, should be considered a major research

445    priority.

446    In conclusion, this work represents the first large scale population genomic study in any

447    *Cryptosporidium* species, inferring the global population structure and evolutionary history of *C.*

448    *hominis*. We propose recognition of two distinct subspecies, *C. h. hominis* (clade 1) and *C. h.*

449    *aquapotentis* (clade 2, *gp60* subtype IbA10G2), with distinct demographic histories that have

450    diverged circa 500 years ago. Although the subspecies differ in their global distribution, their gene

451    pools are not completely isolated, and rare genetic exchanges have occurred in the recent past. We

452    contend that many of the genes, CHUDEA2_430 and CHUDEA6_5270 in particular, in these

453    introgression regions are involved in infection, and that their evolution in humans may be driving

454    greater human specificity, virulence and transmissibility. It appears *C. h. aquapotentis* (clade 2, *gp60*

455    subtype IbA10G2) is playing a key role in this process, which is supported by previous observations

456    based on multilocus typing (Li, et al. 2013). This illustrates how human-mediated gene flow is

457    involved in parasite evolution and genomic architecture, and how it could affect virulence evolution.

458    Also, it shows that the GIPVs that result from population admixture in an anthroponotic species are

459    under selection and involved in evolutionary arms race.

460

461     **Methods**

462     **Parasite isolates**

463     The *C. hominis* isolates newly sequenced for this study (n = 34) were archived stool samples collected

464     at the *Cryptosporidium* Reference Unit in the UK. The species was determined by species-specific

465     real-time PCR targeting the A135 gene (Robinson, et al. 2020) and subtyped by PCR and sequencing

466     of the *gp60* gene (Chalmers, et al. 2019). Supplementary Table 1 provides information about these

467     isolates. Isolates were selected to mainly represent the dominant variant, IbA10G2, as defined by

468     *gp60* sequencing.

469

470     **Processing of faecal samples for whole genome sequencing**

471     Stool samples were processed as previously described (Hadfield, et al. 2015). Briefly, saturated salt-

472     flotation was used to obtain a partially purified suspension of oocysts starting from 1-2 ml of each

473     faecal sample. Oocysts were further purified from the suspension by immunomagnetic separation

474     (IMS), using the Isolate® IMS kit (TCS Biosciences, Botolph Claydon, UK). IMS-purified oocysts

475     were treated with bleach, and washed three times with nuclease-free water by centrifugation at 1,100´

476     g for 5 min. The pellets were suspended in 200 μL of nuclease-free water for DNA extraction.

477

478     **DNA preparation and whole genome sequencing**

479     Genomic DNA was extracted from purified *Cryptosporidium* oocysts by first performing eight cycles

480     of freezing in liquid nitrogen for 1 min and thawing at 95°C for 1 min, and then using the QIAamp

481     DNA extraction kit (Qiagen, Manchester, UK) according to the manufacturer's instructions. The

482     genomic DNA was eluted in 50 μL nuclease-free water, and the concentration measured using the

483     Qubit dsDNA HS Assay Kit with the Qubit 1.0 fluorometer (Invitrogen, Paisley, UK), according to

484     the manufacturer's instructions.

485

486          Whole genome amplification (WGA) was performed using the Repli-g Midi kit (Qiagen,

487     Milan, Italy), according to the manufacturer's instructions. Briefly, 5 μL of genomic DNA (containing

488     1-10 ng of DNA) were mixed with 5 μL of denaturing solution and incubated at room temperature for

489     3 min. Next, 10 μL of stop solution were added to stabilise denatured DNA fragments. The reaction

490     mixture was completed with 29 μL of buffer and 1 μL of phi29 polymerase, and allowed to proceed

491     for 16 hours at 30°C. The reaction was stopped by heating at 63°C for 5 minutes. WGA products were

492     visualised by electrophoresis on a 0.7% agarose gel, purified and quantified by Qubit as described

493     above.

494

495          For Next Generation Sequencing (NGS) experiments, about 1 μg of purified WGA product

496     was used to generate Illumina TruSeq 2x 150 bp paired-end libraries (average insert size: 500 bp),

497     which were sequenced on an Illumina HiSeq 4000 platform (Illumina, SanDiego, CA). Library

498     preparation and NGS experiments were performed by a commercial company (GATC, Germany).

499

500     **Whole genome global dataset**

501     To perform a global comparative genomics of *C. hominis*, we supplemented our newly sequenced

502     genome dataset by downloading all available published *C. hominis* genome sequences on till date

503     (25[th] July 2021), from the sequence read archive (SRA) of NCBI and from the EMBL's European

504     Nucleotide Archive (ENA) (see Supplementary Table 1). Collectively, these data represented 114

505     genome sequences of locally acquired infections from 16 countries across five continents.

506

507     **Data pre-processing and variant calling**

508     Raw reads of the 114 *C. hominis* isolates were trimmed to remove adapter sequences and filtered for

509     low-quality bases using Trimmomatic v.0.36 (Bolger, et al. 2014). The filtered reads were aligned to *C.*

510     *hominis* UdeA01 reference genome (Heiges, et al. 2006; Isaza, et al. 2015) using the maximal exact

511     matches (MEM) algorithm implemented in Burrows-Wheeler Alignment (BWA) tool v.0.7 (Li and

512     Durbin 2009) with default settings. PCR duplicates were then marked using Picard MarkDuplicates

513     (https://broadinstitute.github.io/picard/) followed by Genome Analysis Toolkit's (GATK) indel

514     realignment and base quality score recalibration (BQSR) using default parameters (McKenna, et al.

515     2010). Sequence variants (SNPs) were called from the aligned reads of each isolate using the

516     HaplotypeCaller method in the GATK v3.7.0 (McKenna, et al. 2010) as per GATK's best practices

517     pipeline (Van der Auwera, et al. 2013). SNPs were removed if quality depth (QD) < 2.0, Fisher strand

518     (FS) > 60.0, mapping quality (MQ) < 40.0, mapping quality rank sum test (MQRankSum) < -12.5, read

519     position rank sum test (ReadPosRankSum) < -8.0, Strand odds ratio (SOR) > 4.0. All identified SNPs

520     were combined in one file and each isolate genotyped using the GenotypeGVCFs tool (GATK v3.7.0)

521     (McKenna, et al. 2010). To maximise the quality, SNPs were further filtered based on the following

522     criteria and included in the downstream process: bi-allelic SNPs, quality score > 30, allele depth (AD)

523     > 5, MAF > 0.05 and missing ratio < 0.5. Each of the 114 whole genome sequences assessed here had

524     > 80% coverage of the *C. hominis* reference genome to at least the 5-fold depth. Each of the 114 whole

525     genome sequences assessed here had at least ~80% coverage of the *C. hominis* reference genome to at

526     least the 5-fold depth where 103/114 has > 80% genome coverage and at least 10X coverage. Mean

527     coverage of all isolates is 158X (Quartile1 = 117X, Quartile3 = 229X)] (Supplementary Table 1).

528

529     **Population genetic structure based on whole genome SNPs**

530     The filtered bi-allelic SNPs were used for population structure, phylogenetic and clustering analyses.

531     Multiplicity of infections in each sample were estimated using estMOI (Assefa, et al. 2014) and

532     MOIMIX (https://github.com/bahlolab/moimix). MOIMIX calculates *Fws* statistic (Manske, et al.

533     2012), a fixation index that is used to assess within-host genetic differentiation. An isolate with single

534     infection is expected to have *Fws* 0.95 - 1.00. The R package SNPRelate v.1.18 (Zheng, et al. 2012)

535     was used for principal-component analysis (PCA) analysis. seqVCF2GDS function in SNPRelate R

536     package is used to first convert VCF file into genomic data structure (GDS) file format to store SNP

537     genotypes in an array-oriented matrix format. A genetic covariance matrix is

538     then calculated from genotypes using SNPRelate's function snpgdsPCA, along with the correlation

539     coefficients between samples and genotypes for each SNP. A maximum likelihood phylogenetic tree

540     was constructed by IQ-TREE (Nguyen, et al. 2014) with 1000 bootstraps and visualised in iTOL v3

541     (Letunic and Bork 2016); the sister species *C. parvum* was used as an outgroup. We also constructed a

542     consensus of $10^7$ trees using DensiTree 2 (Bouckaert and Heled 2014) in BEAST v2 (Bouckaert, et al.

543     2014). BEAST v2 (Bouckaert, et al. 2014) was also used to estimate the divergence time between the

544     populations by using 95% highest posterior density (HPD); and SpeedDate

545     (https://github.com/vanOosterhoutLab/SpeedDate.jl) to estimate the coalescence times between

546     sequences by using 5-95% confidence interval (CI). We used mutation rate of $10^{-8}$ and a generation

547     time of 48h/replication (Nader, et al. 2019) to date the coalescence times between sequences.  A

548     Neighbor-Net algorithm-based network was generated using SplitsTree5 (Huson and Bryant 2006).

549     Genetic structure was analysed by STRUCTURE v2.3 software (Pritchard, et al. 2000) for population

550     number (K) ranging 2 - 10 and plotted by using plotSTR R package

551     (https://github.com/DrewWham/Genetic-Structure-Tools). The optimal population genetic cluster

552     value K was estimated by using CLUMPAK (Kopelman, et al. 2015).

553

554     **Population demographic history and divergence time estimation**

555     We used Bayesian Markov Chain Monte Carlo (MCMC) method implemented in Beast v2 program

556     (Bouckaert, et al. 2014) to estimate the effective population size (*Ne*) of the *C. hominis* population.

557     The nucleotide substitution model of HKY was selected. A strict molecular clock model and a

558     Bayesian skyline coalescent tree prior was used with $10^9$ generations of MCMC chain and 10% burn-

559     ins. Tracer v.1.7 (Rambaut, et al. 2018) was used to assess chain convergence and effective sample

560     size [ESS] > 200 and to construct the demographic history over time; i.e. Bayesian Skyline Plot

561     (BSP). SweeD (Pavlidis, et al. 2013) was used to detect windows of selective sweeps from genome-

562     wide SNP dataset by using composite likelihood ratio (CLR) statistic that identifies signature of site

563     frequency spectrum (SFS), with a grid size of 1000.

564

565         Demographic histories and migration rates between the *C. hominis* populations were

566     estimated by using fastsimcoal2 (Excofffier, et al. 2021) by using mutation rate of $10^{-8}$ and a

567     generation time of 48h/replication (Nader, et al. 2019). We first inferred best parameters and the

568     likelihoods for each of the demographic models – no geneflow, ongoing geneflow, early geneflow,

569 recent geneflow and different geneflow, since the time of divergence (~500 years) by running 100

570 independent iterations with 300,000 coalescent simulations and 60 optimisation cycles. Demographic

571 model with the highest likelihood (log10) was then selected to run parameter estimation with block-

572 bootstrapping of 100 replicates.

573

574 **Linkage, recombination and gene flow analyses**

575 We inferred the rate of decay of linkage disequilibrium (LD) by calculating the squared correlation of

576 the coefficient ($r^2$) between SNPs within 50kb using VCFtools (Danecek, et al. 2011). LD blocks were

577 also determined by calculating pairwise $r^2$ between SNPs within chromosomes of each population.

578 Recombination events were identified using the Recombination Detection Program version 4 (RDP4)

579 (Martin, et al. 2015) using the RDP (Martin and Rybicki 2000), GENECONV (Sawyer 1999),

580 BootScan (Salminen, et al. 1995), MaxChi (Smith 1992) and Chimaera (Posada and Crandall 2001)

581 methods. Events were considered significant if at least three methods predicted their occurrence at a

582 probability values, $p \leq 10^{-5}$. Recombination events with undetermined parental sequences were

583 excluded from further HybridCheck analyses. Statistically significant recombination events were

584 visualised and analysed using HybridCheck (Ward and van Oosterhout 2016) to determine the

585 sequence similarity between the isolates involved in the events. HybridCheck program was also used

586 to calculate the D statistic and estimate the gene flow between the populations.

587

588 **Population genetic and genomic analyses of coding region**

589 Tajima's D, Nucleotide diversity ($\pi$), Dxy and Fst were calculated using the PopGenome R-package

590 (Pfeifer, et al. 2014). Nonsynonymous (Ka) and synonymous (Ks) mutation rates were calculated by

591 using Ka/Ks_Calculator (Zhang, et al. 2006). Protein localisation (extracellular) was predicted using

592 WoLF PSORT (Horton, et al. 2007) and information regarding predicted protein targeting (signalling

593 peptides) genes were obtained from CryptoDB (Heiges, et al. 2006). POPART program was used to

594 generate haplotype networks (Leigh and Bryant 2015). AlphaFold was used to predict the protein

595 structures (Jumper, et al. 2021). Glycosylation sites were predicted by using NetNGlyc 4.0 Server

596 (Gupta and Brunak 2002). Intrinsically disordered region in proteins were predicted using IUPred2A

597 (Mészáros, et al. 2018). All statistical tests and results were performed and plotted in R (version

598 3.6.1).

599

600 **Data access**

601 The raw data generated in this study have been submitted to the NCBI BioProject database

602 (https://www.ncbi.nlm.nih.gov/bioproject/) under accession numbers

603    PRJEB15112, PRJNA610731, PRJNA610732, PRJNA610735, PRJNA610737, PRJNA610738,

604    PRJNA610739, PRJNA610740, PRJNA610741, PRJNA610742, PRJNA610743, PRJNA610744,

605    PRJNA610745, PRJNA610746, PRJNA610747 and PRJNA610748. Reviewer link to deposited data

606    for which the accessions are not yet public, are provided in the following table.

| BioProject ID | Reviewer link |
|---|---|
| PRJNA610731 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610731?reviewer=i5832fbel5senhl9g2ju3rj1n1 |
| PRJNA610732 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610732?reviewer=j2pdubv6o5ks7q8earef6kt38r |
| PRJNA610735 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610735?reviewer=g3v1kq0afu72ff6tcc1o61sj5n |
| PRJNA610737 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610737?reviewer=q4joev2kaeck045i21eq6i9s1n |
| PRJNA610738 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610738?reviewer=r4pf8pchijamuem60ncgk1ln01 |
| PRJNA610739 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610739?reviewer=8rgo5b89r49hd0cl27gpmgg1o7 |
| PRJNA610740 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610740?reviewer=rthk0pe2qsn5kc6rn2kdomgo17 |
| PRJNA610741 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610741?reviewer=njvi5jgljq271l2i5ef6hlo0jl |
| PRJNA610742 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610742?reviewer=jdgqt0q8lkgtnt083sdnkrkj0c |
| PRJNA610743 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610743?reviewer=9qka47natbb858bm5fejj5q78q |
| PRJNA610744 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610744?reviewer=8vvg89inuefhack29eo5jscu9b |
| PRJNA610745 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610745?reviewer=lqam5411l8qbhqephabfeq8mu72 |
| PRJNA610746 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610746?reviewer=3hqil4t2u9phuucr8a5nlvnjk0 |
| PRJNA610747 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610747?reviewer=a9bb6grp0g2vsmvljijsgig17t |
| PRJNA610748 | https://dataview.ncbi.nlm.nih.gov/object/PRJNA610748?reviewer=1epesalvpd5qneoadvro9u2436 |

607

627

**Competing interests**

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

630

**Author contributions**

632  A.R.J., S.M.C., C.V.O. and S.T. conceived the study. A.R.J., S.M.C., C.V.O. and S.T. designed the

633  analyses. S.M.C., R.M.C., G.R., D.E., and K.M.T. were involved in acquisition of data. S.T.

634  performed the bioinformatics associated evolutionary genetic and genomic analyses. A.R.J., S.M.C.,

635  C.V.O. and S.T. wrote the manuscript. All authors read and approved the submission of the

636  manuscript for the publication.

637

**Figure Legends**

**Fig 1.** Global population structure of *C. hominis* isolates illustrating their sub-structing and

640  diversification. **a**. PCA of isolates based on the filtered set of whole genome SNPs, highlighting three

641  clusters of isolates which are predominately based on continents of origin. Isolates were color coded

642  with their continent of origin. Isolates associated with *gp60* subtype IbA10G2 were represented with

643  solid circles while non-IbA10G2 with solid triangles. **b.** Structure plot illustrating population genetic

644  ancestry and the admixed nature of the *C. hominis* isolates. The plot was obtained for an optimum

645  value of K=4. The black arrow (bottom) indicates the highly admixed isolate (UK_UKH4), which

646  includes all four ancestries. **c.** Maximum likelihood based phylogenetic tree. **d.** Splitstree and **e.**

647  Densitree are also demonstrating two major clades. *C. h. hominis* (clade 1) includes isolates

648  associated with other *gp60* subtypes while *Cryptosporidium h. aquapotentis* (clade 2) includes

649  isolates associated with *gp60* subtype IbA10G2.

650

**Fig 2**. Demographic histories and population size and secondary contact between *C. h. hominis* (clade

652  1) and *C. h. aquapotentis (*clade 2). **a**. Bayesian Skyline plots (BSP) depicting change in *Ne (*effective

653  population size) through time, for both the clades. The central dark line and the upper and lower

654  dashed lines on Y-axis are mean estimates and 95% HPD intervals of *Ne*, respectively. X-axis is time

655  in years, running backwards. **b**. Boxplot showing significant difference (two-sided t-test) in Tajima's

656 D values between *C. h. hominis (*clade 1) and *C. h. aquapotentis (*clade 2). **c.** Higher likelihood
657 (log10) for "recent geneflow" model (in red). Comparing likelihood distributions of geneflow models
658 and observed significant difference (one-way ANOVA test, F = 2629761, df = 4, p-value <2e-16).
659 Further, Post-hoc Tukey-HSD test revealed difference in likelihood between all the models (p-value <
660 1e-16). **d**. Graphical representation of demographic history of *C. hominis*, illustrating recent
661 secondary contact and migration rates between the two clades (mean ± SE).

662

663 **Fig 3**. Analyses of recombination and gene flow between *C. h. hominis* (clade 1) and *C. h.*
664 *aquapotentis (*clade 2). **a**. Linkage disequilibrium (LD) decay plot showing rapid decay of linkage
665 between SNPs in *C. h. hominis* (clade 1) compared to *C. h. aquapotentis (*clade 2). **b**. Graphical
666 representation of recombinant breakpoint positions detected by RDP4 program between *C. h. hominis*
667 (clade 1) and *C. h. aquapotentis (*clade 2). **c-d**. HybridCheck plots representing genomic signature of
668 introgression in chromosome 2 and 6, respectively. Analysis for chromosome 1 was excluded due to
669 unknown parental sequences. The plots were generated for random set of triplets that includes
670 recombinant (hybrid), minor (donor) and major (recipient) parental sequence, as detected by RDP4
671 program. Introgressed blocks (recombinant breakpoints) were illustrated with dashed boxes, showing
672 high similarity between the recombinant (*C. h. hominis* hybrid isolates) and minor parent (*C. h.*
673 *aquapotentis* isolates). The top panel illustrates the visualisation of sequence similarity between
674 sequences within the triplet, using RBG colour triangle. The two sequences are coloured same
675 (yellow, purple or turquoise) if they share polymorphism. **e.** Gene flow analyses with ABBA-BABA
676 test, representing D statistics for the random sets of triplets (as used in c-d) along with *C. parvum* as
677 an outgroup. D statistic values close to -1 at all three recombinant events, suggesting geneflow
678 between H1 and H3. **f-g**. Pairwise LD of SNPs in chromosomes 2 and 6 of *C. h. hominis* showing red-
679 blocks of high linkage between SNPs in introgressed events 2-4.

680

681 **Fig 4**. Population genetic analyses of genomic islands of putative virulence genes (GIPVs). **a**.
682 Population genetic and divergence analyses of introgressed regions. X-axis represents genomic
683 positions of eight chromosomes highlighted with different colours. Population divergence (Dxy)
684 between *C. h. hominis (*clade 1) and *C. h. aquapotentis (*clade 2) for each gene were plotted on Y-axis
685 (top panel). Nucleotide diversity (π) for *C. h. hominis* (middle panel) and *C. h. aquapotentis* (bottom
686 panel) for each gene, were also plotted on Y-axis, respectively. The breakpoints of four recombination
687 events (event 1-4) were indicated by grey vertical boxes. Event 1 was un-detected in *C. h.*
688 *aquapotentis*. **b**. Correlation between π and Dxy were plotted to identify polymorphic and potential
689 virulence genes.

690

691    **Fig 5.** Illustrating diversifying selection between *C. hominis* subspecies and host adaptation at

692    CHUDEA6_5270 (hypothetical gene). **a.** Haplotype network analyses illustrating haplotype

693    diversification between *C. h. hominis (*clade 1) and *C. h. aquapotentis (*clade 2). **b**. Pairwise

694    nucleotide divergence shows bimodal distribution, which, theoretically, can be explained both by

695    balancing selection (Lighten, et al. 2017), as well as by genetic introgression. **c.** Comparison of

696    predicted models of protein structure of CHUDEA6_5270 gene between *Cryptosporidium* species and

697    subtypes demonstrates variation towards C.-terminal region. **d**. Introgressed-isolates driving balancing

698    selection at gene CHUDEA6_5270 in *C. h. hominis*. Red line represents balancing selection (positive

699    Tajima's D) in *C. h. hominis* that also includes introgressed-isolates. Blue line represents purifying

700    selection (negative Tajima's D) in *C. h. hominis* after excluding introgressed-isolates.

701

702    **References**

703    Abrahamsen MS, Templeton TJ, Enomoto S, Abrahante JE, Zhu G, Lancto CA, Deng M, Liu C,

704    Widmer G, Tzipori S, et al. 2004. Complete genome sequence of the apicomplexan, *Cryptosporidium*

705    *parvum*. Science 304:441-445.

706

707    Assefa SA, Preston MD, Campino S, Ocholla H, Sutherland CJ, Clark TG. 2014. estMOI: estimating

708    multiplicity of infection using parasite deep sequencing data. Bioinformatics 30:1292-1294.

709

710    Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina sequence data.

711    Bioinformatics 30:2114-2120.

712

713    Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, Suchard MA, Rambaut A, Drummond

714    AJ. 2014. BEAST 2: A Software Platform for Bayesian Evolutionary Analysis. PLOS Computational

715    Biology 10:e1003537.

716

717    Bouckaert RR, Heled J. 2014. DensiTree 2: Seeing Trees Through the Forest. bioRxiv:012401.

718    Bouzid M, Hunter PR, Chalmers RM, Tyler KM. 2013. *Cryptosporidium* pathogenicity and virulence.

719    Clin Microbiol Rev 26:115-134.

720

721    Cacciò SM, Chalmers RM. 2016. Human cryptosporidiosis in Europe. Clin Microbiol Infect 22:471-

722    480.

723

724    Cama VA, Bern C, Roberts J, Cabrera L, Sterling CR, Ortega Y, Gilman RH, Xiao L. 2008.

725    *Cryptosporidium* species and subtypes and clinical manifestations in children, Peru. Emerging

726    infectious diseases 14:1567.

727

728 Carmicheal J, Atri P, Sharma S, Kumar S, Chirravuri Venkata R, Kulkarni P, Salgia R, Ghersi D, Kaur
729 S, Batra SK. 2020. Presence and structure-activity relationship of intrinsically disordered regions across
730 mucins. The FASEB Journal 34:1939-1957.

731

732 Cevallos AM, Bhat N, Verdon R, Hamer DH, Stein B, Tzipori S, Pereira ME, Keusch GT, Ward HD.
733 2000. Mediation of *Cryptosporidium parvum* infection in vitro by mucin-like glycoproteins defined by
734 a neutralizing monoclonal antibody. Infect Immun 68:5167-5175.

735

736 Chalmers RM, Robinson G, Elwin K, Elson R. 2019. Analysis of the *Cryptosporidium* spp. and *gp60*
737 subtypes linked to human outbreaks of cryptosporidiosis in England and Wales, 2009 to 2017. Parasites
738 & Vectors 12:95.

739

740 Chalmers RM, Robinson G, Elwin K, Hadfield SJ, Thomas E, Watkins J, Casemore D, Kay D. 2010.
741 Detection of *Cryptosporidium* species and sources of contamination with *Cryptosporidium hominis*
742 during a waterborne outbreak in north west Wales. Journal of water and health 8:311-325.

743

744 Chatterjee A, Banerjee S, Steffen M, O'Connor RM, Ward HD, Robbins PW, Samuelson J. 2010.
745 Evidence for mucin-like glycoproteins that tether sporozoites of *Cryptosporidium parvum* to the inner
746 surface of the oocyst wall. Eukaryotic cell 9:84-96.

747

748 Corsi GI, Tichkule S, Sannella AR, Vatta P, Asnicar F, Segata N, Jex AR, Oosterhout Cv, Cacciò SM.
749 2021. Evolutionary epidemiology of a zoonosis. bioRxiv:2021.2010.2015.464618.

750

751 Craun GF, Hubbs SA, Frost F, Calderon RL, Via SH. 1998. Waterborne outbreaks of cryptosporidiosis.
752 Journal AWWA 90:81-91.

753

754 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE, Lunter G, Marth
755 GT, Sherry ST, et al. 2011. The variant call format and VCFtools. Bioinformatics 27:2156-2158.

756

757 Drummond AJ, Rambaut A, Shapiro B, Pybus OG. 2005. Bayesian coalescent inference of past
758 population dynamics from molecular sequences. Mol Biol Evol 22:1185-1192.

759

760 Efstratiou A, Ongerth JE, Karanis P. 2017. Waterborne transmission of protozoan parasites: review of
761 worldwide outbreaks-an update 2011–2016. Water research 114:14-22.

762

763    Excofffier L, Marchi N, Marques DA, Matthey-Doret R, Gouy A, Sousa VC. 2021. fastsimcoal2:
764    demographic inference under complex evolutionary scenarios. Bioinformatics.
765

766    Feng Y, Ryan UM, Xiao L. 2018. Genetic Diversity and Population Structure of *Cryptosporidium*.
767    Trends Parasitol.
768

769    Feng Y, Tiao N, Li N, Hlavsa M, Xiao L. 2014. Multilocus sequence typing of an emerging
770    *Cryptosporidium hominis* subtype in the United States. J Clin Microbiol 52:524-530.
771

772    Guo Y, Tang K, Rowe LA, Li N, Roellig DM, Knipe K, Frace M, Yang C, Feng Y, Xiao L. 2015.
773    Comparative genomic analysis reveals occurrence of genetic recombination in virulent
774    *Cryptosporidium hominis* subtypes and telomeric gene duplications in *Cryptosporidium parvum*. BMC
775    Genomics 16:320.
776

777    Gupta R, Brunak S. 2002. Prediction of glycosylation across the human proteome and the correlation
778    to protein function. Pac Symp Biocomput:310-322.
779

780    Hadfield SJ, Pachebat JA, Swain MT, Robinson G, Cameron SJ, Alexander J, Hegarty MJ, Elwin K,
781    Chalmers RM. 2015. Generation of whole genome sequences of new *Cryptosporidium hominis* and
782    *Cryptosporidium parvum* isolates directly from stool samples. BMC Genomics 16:650.
783

784    Heiges M, Wang H, Robinson E, Aurrecoechea C, Gao X, Kaluskar N, Rhodes P, Wang S, He CZ, Su
785    Y, et al. 2006. CryptoDB: a *Cryptosporidium* bioinformatics resource update. Nucleic Acids Res
786    34:D419-422.
787

788    Herges GR, Widmer G, Clark ME, Khan E, Giddings CW, Brewer M, McEvoy JM. 2012. Evidence
789    that *Cryptosporidium parvum* populations are panmictic and unstructured in the Upper Midwest of the
790    United States. Appl Environ Microbiol 78:8096-8101.
791

792    Horton P, Park K-J, Obayashi T, Fujita N, Harada H, Adams-Collier CJ, Nakai K. 2007. WoLF PSORT:
793    protein localization predictor. Nucleic Acids Res 35:W585-W587.
794

795    Huson DH, Bryant D. 2006. Application of phylogenetic networks in evolutionary studies. Mol Biol
796    Evol 23:254-267.
797

798    Isaza JP, Galván AL, Polanco V, Huang B, Matveyev AV, Serrano MG, Manque P, Buck GA, Alzate
799    JF. 2015. Revisiting the reference genomes of human pathogenic *Cryptosporidium* species:
800    reannotation of *C. parvum* Iowa and a new *C. hominis* reference. Sci Rep 5:16324.
801
802    Jex AR, Gasser RB. 2010. Genetic richness and diversity in *Cryptosporidium hominis* and *C. parvum*
803    reveals major knowledge gaps and a need for the application of "next generation" technologies--
804    research review. Biotechnol Adv 28:17-26.
805
806    Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates R,
807    Žídek A, Potapenko A, et al. 2021. Highly accurate protein structure prediction with AlphaFold. Nature.
808
809    Khalil IA, Troeger C, Rao PC, Blacker BF, Brown A, Brewer TG, Colombara DV, De Hostos EL,
810    Engmann C, Guerrant RL, et al. 2018. Morbidity, mortality, and long-term consequences associated
811    with diarrhoea from *Cryptosporidium* infection in children younger than 5 years: a meta-analyses study.
812    Lancet Glob Health 6:e758-e768.
813
814    King P, Robinson G, Elwin K, Tyler KM, Hunter PR, Chalmers RM. 2017. Prevalence and
815    epidemiology of human *Cryptosporidium parvum* IIc infections in England and Wales. The Lancet
816    389:S56.
817
818    King P, Tyler KM, Hunter PR. 2019. Anthroponotic transmission of *Cryptosporidium parvum*
819    predominates in countries with poorer sanitation: a systematic review and meta-analysis. Parasites &
820    Vectors 12:16-16.
821
822    Kopelman NM, Mayzel J, Jakobsson M, Rosenberg NA, Mayrose I. 2015. Clumpak: a program for
823    identifying clustering modes and packaging population structure inferences across K. Molecular
824    ecology resources 15:1179-1191.
825
826    Kotloff KL, Nataro JP, Blackwelder WC, Nasrin D, Farag TH, Panchalingam S, Wu Y, Sow SO, Sur
827    D, Breiman RF, et al. 2013. Burden and aetiology of diarrhoeal disease in infants and young children
828    in developing countries (the Global Enteric Multicenter Study, GEMS): a prospective, case-control
829    study. Lancet 382:209-222.
830
831    Leigh JW, Bryant D. 2015. popart: full-feature software for haplotype network construction. Methods
832    in Ecology and Evolution 6:1110-1116.
833

834    Letunic I, Bork P. 2016. Interactive tree of life (iTOL) v3: an online tool for the display and annotation
835    of phylogenetic and other trees. Nucleic Acids Res 44:W242-245.
836

837    Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
838    Bioinformatics 25:1754-1760.
839

840    Li N, Xiao L, Cama VA, Ortega Y, Gilman RH, Guo M, Feng Y. 2013. Genetic recombination and
841    *Cryptosporidium hominis* virulent subtype IbA10G2. Emerg Infect Dis 19:1573-1582.
842

843    Lighten J, Papadopulos AST, Mohammed RS, Ward BJ, G. Paterson I, Baillie L, Bradbury IR, Hendry
844    AP, Bentzen P, van Oosterhout C. 2017. Evolutionary genetics of immunological supertypes reveals
845    two faces of the Red Queen. Nature Communications 8:1294.
846

847    Lin B, Qing X, Liao J, Zhuo K. 2020. Role of Protein Glycosylation in Host-Pathogen Interaction. Cells
848    9:1022.
849

850    Lippuner C, Ramakrishnan C, Basso WU, Schmid MW, Okoniewski M, Smith NC, Hässig M, Deplazes
851    P, Hehl AB. 2018. RNA-Seq analysis during the life cycle of *Cryptosporidium parvum* reveals
852    significant differential gene expression between proliferating stages in the intestine and infectious
853    sporozoites. Int J Parasitol 48:413-422.
854

855    Manske M, Miotto O, Campino S, Auburn S, Almagro-Garcia J, Maslen G, O'Brien J, Djimde A,
856    Doumbo O, Zongo I, et al. 2012. Analysis of *Plasmodium falciparum* diversity in natural infections by
857    deep sequencing. Nature 487:375-379.
858

859    Martin D, Rybicki E. 2000. RDP: detection of recombination amongst aligned sequences.
860    Bioinformatics 16:562-563.
861

862    Martin DP, Murrell B, Golden M, Khoosal A, Muhire B. 2015. RDP4: Detection and analysis of
863    recombination patterns in virus genomes. Virus Evolution 1.
864

865    McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D,
866    Gabriel S, Daly M. 2010. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-
867    generation DNA sequencing data. Genome research 20:1297-1303.
868

869 McKerr C, Chalmers RM, Elwin K, Jones H, Vivancos R, O'Brien SJ, Christley RM. 2021. Cross-
870 Sectional Household Transmission Study of *Cryptosporidium* Shows that *C. hominis* Infections are A
871 Key Risk Factor for Spread.
872
873 Mészáros B, Erdős G, Dosztányi Z. 2018. IUPred2A: context-dependent prediction of protein disorder
874 as a function of redox state and protein binding. Nucleic Acids Res 46:W329-W337.
875
876 Morrison LJ, Mallon ME, Smith HV, MacLeod A, Xiao L, Tait A. 2008. The population structure of
877 the *Cryptosporidium parvum* population in Scotland: a complex picture. Infect Genet Evol 8:121-129.
878
879 Nader JL, Mathers TC, Ward BJ, Pachebat JA, Swain MT, Robinson G, Chalmers RM, Hunter PR, van
880 Oosterhout C, Tyler KM. 2019. Evolutionary genomics of anthroponosis in *Cryptosporidium*. Nature
881 Microbiology.
882
883 Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2014. IQ-TREE: A Fast and Effective Stochastic
884 Algorithm for Estimating Maximum-Likelihood Phylogenies. Mol Biol Evol 32:268-274.
885
886 Nichols GL, Chalmers RM, Hadfield SJ. 2014. Molecular epidemiology of human cryptosporidiosis.
887 In. *Cryptosporidium*: parasite and disease: Springer. p. 81-147.
888
889 O'Connor RM, Burns PB, Ha-Ngoc T, Scarpato K, Khan W, Kang G, Ward H. 2009. Polymorphic
890 mucin antigens CpMuc4 and CpMuc5 are integral to *Cryptosporidium parvum* infection in vitro.
891 Eukaryotic cell 8:461-469.
892
893 Pavlidis P, Živkovic D, Stamatakis A, Alachiotis N. 2013. SweeD: likelihood-based detection of
894 selective sweeps in thousands of genomes. Mol Biol Evol 30:2224-2234.
895
896 Perez-Vilar J, Hill RL. 1999. The Structure and Assembly of Secreted Mucins* 210. Journal of
897 Biological Chemistry 274:31751-31754.
898
899 Pfeifer B, Wittelsbürger U, Ramos-Onsins SE, Lercher MJ. 2014. PopGenome: an efficient Swiss army
900 knife for population genomic analyses in R. Mol Biol Evol 31:1929-1936.
901
902 Pond SL, Frost SD. 2005. Datamonkey: rapid detection of selective pressure on individual sites of codon
903 alignments. Bioinformatics 21:2531-2533.
904

905      Posada D, Crandall KA. 2001. Evaluation of methods for detecting recombination from DNA
906      sequences: Computer simulations. Proceedings of the National Academy of Sciences 98:13757.
907

908      Pritchard JK, Stephens M, Donnelly P. 2000. Inference of population structure using multilocus
909      genotype data. Genetics 155:945-959.
910

911      Putignani L, Menichella D. 2010. Global distribution, public health and clinical impact of the protozoan
912      pathogen *Cryptosporidium*. Interdiscip Perspect Infect Dis 2010.
913

914      Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. 2018. Posterior Summarization in Bayesian
915      Phylogenetics Using Tracer 1.7. Syst Biol 67:901-904.
916

917      Razakandrainibe R, Diawara EHI, Costa D, Le Goff L, Lemeteil D, Ballet JJ, Gargala G, Favennec L.
918      2018. Common occurrence of *Cryptosporidium hominis* in asymptomatic and symptomatic calves in
919      France. PLoS Negl Trop Dis 12:e0006355.
920

921      Robinson G, Elwin K, Chalmers RM. 2020. *Cryptosporidium* Diagnostic Assays: Molecular Detection.
922      Methods Mol Biol 2052:11-22.
923

924      Ryan U, Zahedi A, Paparini A. 2016. *Cryptosporidium* in humans and animals—a one health approach
925      to prophylaxis. Parasite Immunol 38:535-547.
926

927      Salminen MO, CARR JK, BURKE DS, McCUTCHAN FE. 1995. Identification of breakpoints in
928      intergenotypic recombinants of HIV type 1 by bootscanning. AIDS research and human retroviruses
929      11:1423-1425.
930

931      Santin M. 2020. *Cryptosporidium* and *Giardia* in Ruminants. Vet Clin North Am Food Anim Pract
932      36:223-238.
933

934      Sawyer S. 1999. GENECONV: a computer package for the statistical detection of gene conversion.
935      http://www.math.wustl.edu/~ sawyer.
936

937      Segura R, Prim N, Montemayor M, Valls ME, Muñoz C. 2015. Predominant virulent IbA10G2 subtype
938      of *Cryptosporidium hominis* in human isolates in Barcelona: a five-year study. PLoS One 10:e0121753.
939

940      Smith JM. 1992. Analyzing the mosaic structure of genes. J Mol Evol 34:126-129.
941

942 Tichkule S, Jex AR, van Oosterhout C, Sannella AR, Krumkamp R, Aldrich C, Maiga-Ascofare O,
943 Dekker D, Lamshöft M, Mbwana J, et al. 2021. Comparative genomics revealed adaptive admixture in
944 *Cryptosporidium hominis* in Africa. Microb Genom 7.

946 Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, Jordan T,
947 Shakir K, Roazen D, Thibault J, et al. 2013. From FastQ data to high confidence variant calls: the
948 Genome Analysis Toolkit best practices pipeline. Curr Protoc Bioinformatics 43:11.10.11-33.

950 Van Oosterhout C. 2021. Mitigating the threat of emerging infectious diseases; a coevolutionary
951 perspective. Virulence 12:1288-1295.

953 Wang R, Zhang L, Axen C, Bjorkman C, Jian F, Amer S, Liu A, Feng Y, Li G, Lv C, et al. 2014.
954 *Cryptosporidium parvum* IId family: clonal population and dispersal from Western Asia to other
955 geographical regions. Sci Rep 4:4208.

957 Ward BJ, van Oosterhout C. 2016. HYBRIDCHECK: software for the rapid detection, visualization
958 and dating of recombinant regions in genome sequence data. Mol Ecol Resour 16:534-539.

960 Widerström M, Schönning C, Lilja M, Lebbad M, Ljung T, Allestam G, Ferm M, Björkholm B, Hansen
961 A, Hiltula J. 2014. Large outbreak of *Cryptosporidium hominis* infection transmitted through the public
962 water supply, Sweden. Emerging infectious diseases 20:581.

964 Yang X, Guo Y, Xiao L, Feng Y. 2021. Molecular epidemiology of human cryptosporidiosis in low-
965 and middle-income countries. Clinical Microbiology Reviews 34:e00087-00019.

967 Zhang Z, Li J, Zhao XQ, Wang J, Wong GK, Yu J. 2006. KaKs_Calculator: calculating Ka and Ks
968 through model selection and model averaging. Genomics Proteomics Bioinformatics 4:259-263.

970 Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS. 2012. A high-performance computing
971 toolset for relatedness and principal component analysis of SNP data. Bioinformatics 28:3326-3328.

973 Zhou L, Singh A, Jiang J, Xiao L. 2003. Molecular surveillance of *Cryptosporidium* spp. in raw
974 wastewater in Milwaukee: implications for understanding outbreak occurrence and transmission
975 dynamics. J Clin Microbiol 41:5254-5257.

**Figures**

**Figure 1.**

991        **Figure 2.**



992

993

994

995

996

997

998    **Figure 3.**



999

1000

1001

1002

1003

1004

1005

1006

1007    **Figure 4.**

1012    **Figure 5.**



1013

1014

1015