

1 **Title: Re-annotation of the *Theileria parva* genome refines 53% of the proteome and**
2 **uncovers essential components of N-glycosylation, a conserved pathway in many**
3 **organisms**

4 Kyle Tretina¹, Roger Pelle², Joshua Orvis¹, Hanzel T. Gotia¹, Olukemi O. Ifeonu¹, Priti
5 Kumari¹, Nicholas C. Palmateer¹, Shaikh B.A. Iqbal¹, Lindsay Fry^{3,4}, Vishvanath M.
6 Nene⁵, Claudia Daubenberger⁶, Richard P. Bishop³, Joana C. Silva^{1,7*}

7 **Affiliations:**

8 ¹Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore,
9 MD, USA

10 ²Biosciences eastern and central Africa-International Livestock Research Institute,
11 Nairobi, Kenya

12 ³Animal Disease Research Unit, Agricultural Research Service, USDA, Pullman, WA
13 99164-7030, USA

14 ⁴Department of Veterinary Microbiology & Pathology, Washington State University
15 Pullman, WA 99164-7040, USA

16 ⁵International Livestock Research Institute, Nairobi, Kenya

17 ⁶Swiss Tropical and Public Health Institute, Basel, Switzerland & University of Basel,
18 Basel, Switzerland

19 ⁷Department of Microbiology and Immunology, University of Maryland School of
20 Medicine, Baltimore, MD, USA

21 *Corresponding author

22

23

24 **Abstract (<350 words)**

25 Background: Genome annotation remains a significant challenge because of limitations in
26 the quality and quantity of the data being used to inform the location and function of
27 protein-coding genes and, when RNA data are used, the underlying biological complexity
28 of the processes involved in gene expression. However, comprehensive and descriptive
29 genome annotations are valuable resources for answering biological research questions
30 and discovering novel chemotherapeutic targets for disease treatment.

31 Results: Here, we apply our recently published RNAseq dataset derived from the schizont
32 life-cycle stage of the apicomplexan parasite *Theileria parva*, which causes a devastating
33 livestock disease in sub-Saharan Africa, to update structural and functional annotations
34 across the entire nuclear genome.

35 Conclusions: The process of re-annotation led to novel insights into the organization and
36 expression profiles of protein-coding sequences in this parasite, and uncovered a minimal
37 N-glycosylation pathway that changes our current understanding of the evolution of this
38 post-translation modification in apicomplexan parasites.

39

40 **Key Words:** *Theileria*, East coast fever, Genome, Re-annotation, N-glycosylation

41

42 **Background**

43 East Coast fever (ECF) in eastern, central, and southern Africa causes an
44 estimated loss of over 1 million heads of cattle yearly, with an annual economic loss that
45 surpasses \$300 million USD, impacting mainly smallholder farmers [1]. Cattle are the
46 most valuable possession of smallholder farmers in this region, as they are a source of

47 milk, meat and hides, provide manure and traction in mixed crop-livestock systems, and
48 revenue derived from livestock pays for school fees and dowries [2, 3]. ECF is a tick-
49 transmitted disease caused by the apicomplexan parasite *Theileria parva*. Lymphocytes
50 infected with *T. parva* proliferate in the regional lymph node draining the tick bite site,
51 and then metastasize into various lymphoid and non-lymphoid organs, and induce a
52 severe inflammatory reaction that leads to respiratory failure and death of susceptible
53 cattle, which typically die within three to four weeks of infection [4-7]. *T. parva* control
54 is vital to food security in this region of the world, which is plagued by a range of other
55 infectious diseases of humans and their livestock.

56 Efficacious and affordable chemotherapeutics and vaccines are essential tools in
57 the effective control of infectious disease agents [8, 9]. A reliable structural annotation of
58 the genome, consisting at minimum of the correct location of all protein-coding
59 sequences (CDSs), enables the identification, prioritization and experimental screening of
60 potential vaccine and drug targets [10-12]. The accurate identification of the complete
61 proteome can greatly enhance microbiological studies, and reveals metabolic processes
62 unique to pathogens [13]. In turn, a better understanding of the biology of *T. parva*
63 transmission, colonization and pathogenesis may ultimately reveal novel targets for
64 pathogen control [14]. Currently, much like for other apicomplexan parasites [15, 16],
65 knowledge on the functional role of genomic sequences outside of *T. parva* CDSs is
66 sparse, and many gene models containing only CDSs are supported by little or no
67 experimental evidence. RNAseq data, generated through deep sequencing of cDNA using
68 next generation sequencing technologies, can provide an extraordinary level of insight
69 into gene structure and regulation [12, 17]. Here, we used the first high-coverage

70 RNAseq data for this species [18] to improve existing gene models through the
71 identification of start and stop codons, primary intron splice sites and untranslated
72 regions (UTRs). This new gene model annotation brought to light several new insights
73 into gene expression in this gene-dense eukaryote, and led to the discovery of several
74 new prospective chemotherapeutic targets for treating ECF.

75

76 **Results**

77 *The annotation of the Theileria parva genome is significantly improved, revealing a*
78 *higher gene density than previously thought*

79 The nuclear genome of the reference *T. parva* Muguga isolate consists of four
80 linear chromosomes which are currently assembled into eight contigs (Supplementary
81 Table S1, Additional File 1): chromosomes 1 and 2 are assembled into a single contig
82 each, chromosome 3 is in four contigs and chromosome 4 in two [19]. The new genome
83 annotation was based on this assembly and on extensive RNAseq data (Supplementary
84 Figure S1, Additional File 1). We performed a comprehensive revision of the entire *T.*
85 *parva* genome annotation, including automated structural annotation and a double-pass
86 manual curation of each locus (see Methods).

87 The re-annotation process resulted in the discovery of 128 new genes, 274
88 adjacent gene models were merged, 157 gene models were split, and 38 genes were
89 replaced by new genes encoded in the reverse orientation (Figure 1). In addition, exons
90 boundaries have been corrected in over a thousand genes. Overall, 83% of all nuclear
91 genes in the original annotation were altered in some way, with changes made on every

92 contig. This resulted in significant alterations to the predicted proteome, with 53% of the
93 nuclear proteins in the original annotation having altered amino acid sequences in the
94 new annotation, a remarkable ~50 bp increase in average CDSs length, a reduction of the
95 average length of intergenic regions by close to 100 bp and the assignment of an
96 additional 200,000 base pairs (or 2.4% of the genome), previously classified as intergenic
97 or intronic sequences, to the proteome. This results in a genome that is denser than
98 previously thought, with an overall increase in the coding fraction of the genome from
99 68% to 71%, more closely resembling *T. annulata* Ankara, which has a coding fraction of
100 72.9% (Supplementary Table S2, Additional File 1). In fact, *T. parva* has the densest
101 genome out of the indicated genomes investigated, with one protein-coding gene every
102 ~2100 bp (Supplementary Table S2, Additional File 1).

103 Several lines of evidence suggest that this annotation represents a very significant
104 improvement of the *T. parva* proteome relative to the original annotation. First, there was
105 an increase in the proportion of proteins with at least one PFAM domain in the new
106 proteome compared to the original proteome, implying that the new annotation captures
107 functional elements that were previously missed (Figure 2a). Given the close
108 evolutionary relationship and near complete synteny between *T. parva* and *T. annulata*
109 [20], their respective proteomes are expected to be very similar. Indeed, a comparison of
110 the two predicted proteomes results in 52 additional reciprocal best hits and protein
111 length differences between orthologs in *T. parva* and *T. annulata* also decreased
112 significantly (Figure 2b). It is likely that some of the most significant differences between
113 the *T. parva* and *T. annulata* proteomes, in particular the 25% fewer protein-coding genes
114 and much longer CDSs in the latter, represent annotation errors in the *T. annulata*

115 genome that will be corrected upon revision with more recently accumulated evidence.

116 The total number of non-canonical splice sites in the genome increased from 0.15% to

117 0.36% of all introns, but the sequence diversity of non-canonical splice sites decreased

118 from eight non-canonical splice donor and acceptor site combinations to only a single

119 splice site pair – GC/AG donor and acceptor dinucleotides, recognized by the U2-type

120 spliceosome [21] (Figure 2c). The new annotation is also considerably more consistent

121 with the RNAseq data, with a larger number of introns, a higher proportion of which is

122 supported by at least one RNAseq read (Figure 2d). A total of 118 introns from the

123 original genome annotation have been removed, due to contradicting RNAseq evidence.

124 The tremendous power of RNAseq to inform on gene and isoform structure in this

125 species revealed a significant amount of transcriptome diversity and complexity. First,

126 the proportion of loci, defined here as a continuous genomic region encoding the length

127 of a CDS, intervening introns, and flanking UTRs, that appear to overlap an adjacent

128 locus increased from 2% to 10% in the new annotation. In many of these instances, read

129 coverage, coding potential, and other evidence support the presence of adjacent genes

130 with overlapping UTRs. In 125 cases, the overlap includes not only UTRs but also CDSs.

131 Secondly, there are many instances of overlapping loci in which the respective CDSs are

132 encoded in the same strand; in these cases, no UTRs were defined in the intervening

133 intergenic region, since their exact boundaries could not be determined. Finally, during

134 manual curation, we observed many instances of potential alternative splicing, the

135 clearest of which were the cases of well-supported introns where RNAseq coverage was

136 nevertheless significantly higher than zero. In these cases, only the most prevalent

137 isoform was annotated (Figure 1f). Finally, despite its power, RNAseq evidence is not

138 sufficient to resolve the structure of all loci; when the evidence did not clearly favor one
139 gene model over another, the gene model in the original annotation was maintained by
140 default. Interestingly, the vast majority of the genes appear to have only one or,
141 sometimes, two most prevalent isoforms, as has been proposed for *Plasmodium* [22],
142 although this was not defined quantitatively here. The median length of the annotated
143 mRNA reported here is ~1,500 bp, and the maximum length >15,000 bp (Supplementary
144 Figure S2, Additional File 1).

145

146 *Most genes are transcribed during the schizont stage of the Theileria parva life-cycle, and*
147 *antisense transcription is widespread*

148 We sequenced cDNA generated from polyA-enriched total RNA collected from a
149 *T. parva*-infected, schizont-transformed bovine cell line (see Methods section). A total of
150 8.3×10^7 paired-end reads were obtained with an Illumina HiSeq 2000 platform, 70.04%
151 of which mapped to the *T. parva* reference genome (Supplementary Table S1, Additional
152 File 1). RNAseq provided a complete and quantitative view of transcription revealing that
153 most of the genome of this parasite is transcribed during the schizont stage of its life
154 cycle (Supplementary Figures S1, S3, Additional File 1). We found that 4011 of all 4054
155 (98%) predicted protein-coding parasite genes are transcribed at the schizont stage, and
156 12,172 of all 12,296 introns are supported by RNAseq reads (Figure 2d). We found
157 evidence of expression for almost all of the known humoral and cellular immunity
158 antigens (Supplementary Table S3, Additional File 1). In fact, Tp9, one of those antigens,
159 is among the 15 most highly expressed genes in our dataset (Supplementary Table S4,

160 Additional File 1). Interestingly, its ortholog in *T. annulata* has been hypothesized to
161 contribute to schizont-induced host cell transformation [23].

162 As has recently been suggested from *in silico* analyses [18], transcription in *T.*
163 *parva* occurs from diverse kinds of promoters, with many instances of adjacent loci
164 overlapping on the same or opposite strands. In fact, of the 4,085 predicted protein-
165 coding nuclear genes, only 74 had an estimated reads per kilobase of transcript per
166 million reads (RPKM) of zero and an additional 154 had RPKM<1. Interestingly, of the
167 74 genes with an RPKM of zero, most are hypothetical, with no predicted functional
168 annotation, and without any high-confidence orthologs (Supplementary Table S5,
169 Additional File 1). Since tRNAs are not polyadenylated, they were not found in our
170 RNAseq dataset (Materials and Methods). Annotated protein-coding genes lacking
171 RNAseq evidence are mostly orthologs of *Plasmodium falciparum* apicoplast proteins
172 with mid blood stage expression [24, 25], *T. parva* repeat (*Tpr*) family proteins, or
173 DUF529 domain-containing proteins (Supplementary Table S5, Additional File 1). These
174 data are consistent with a study published in 2005, which used MPSS to estimate
175 expression levels of *T. parva* genes in the schizont stage of the parasite [26], as well as a
176 more recent study comparing gene expression between the schizont and the
177 sporozoite/sporoblast stages [27]. The expression levels in the sense strand for each gene,
178 as quantified by RPKM, when log-transformed, followed a unimodal distribution similar
179 to a normal distribution (Figure 3a).

180

181 *T. parva* multi-gene families show variable expression levels

182 Large gene families are known to play a role in the pathogenesis of protozoan
183 infections, perhaps the most well-known being the *var* gene family in *P. falciparum*.
184 These genes encode proteins that are essential for the sequestration of infected red blood
185 cells, a critical biological feature determining severe malaria pathology of *P. falciparum*
186 [28]. Using the OrthoMCL algorithm as described previously [19], we clustered paralogs
187 in this genome, identifying changes in the size of several of the largest *T. parva* gene
188 families (Supplementary Table S6, Additional File 1), and finding variable patterns in
189 their levels of expression (Supplementary Figure S4, Additional File 1). The roles of
190 most of these gene families are not known. For example, the *Tpr* (*T. parva* repeat) gene
191 family has been suggested to be rapidly evolving and expressed as protein in the
192 piroplasm stages [19]. This is consistent with our findings, which show *Tpr* genes not to
193 be highly expressed in the schizont (Supplementary Figure S4, Additional File 1) or the
194 sporoblast (Supplementary Figure S5, Additional File 1) stages [27, 29]. Interestingly, in
195 that same dataset, we find a significant up-regulation of subtelomeric variable secreted
196 protein gene (SVSP) family genes in the sporozoite stages relative to both the sporoblast
197 and schizont stages, suggesting that they may be important for invasion or the
198 establishment of infection in the vertebrate host (Supplementary Figure S5, Additional
199 File 1) [30].

200 This new *T. parva* genome annotation not only improved our resolution of the
201 gene models of multi-gene family members and other transformation factors
202 (Supplementary Figure S6, Additional File 1) [31], but also uncovered 128 genes that
203 were not present in the original annotation.

204

205 *A mechanism of core N-glycosylation is now predicted in T. parva*

206 Among the 128 newly identified genes, one was annotated as a potential Alg14
207 ortholog, an important part of a glycosyltransferase complex in many organisms that add
208 a N-acetylglucosamine (GlcNAc) to the N-glycan precursor. N-glycosylation is an
209 important type of protein post-translation modification, during which a sugar is linked to
210 the nitrogen of specific amino acid residues, a process that occurs in the membrane of the
211 endoplasmic reticulum and is critical for both the structure and function of many
212 eukaryotic proteins. N-glycosylation is a ubiquitous protein modification process, but the
213 glycans being transferred differ among the domains of life [32]. However, in
214 apicomplexan parasites that infect red blood cells, there appears to be a selection against
215 long N-glycan chains [33]. *Theileria* parasites were previously believed to not add N-
216 acetylglucosamine to their glycan precursors, since sequence similarity searches did not
217 identify the necessary enzymes. While this previous study did not discover any Alg
218 enzymes, we find that *T. parva* has Alg7 (*TpAlg7*; TpMuguga_01g00118), Alg13
219 (*TpAlg13*; TpMuguga_02g00515), and Alg14 (*TpAlg14*; TpMuguga_01g02045)
220 homologs, which show differential mRNA-level expression between the sporozoite and
221 schizont life cycle stages (Supplementary Figure S7, Additional File 1). In fact, the
222 structure of each of these *Theileria* proteins can be predicted *ab initio* with high
223 confidence (Supplementary Table S7, Additional File 1) and have predicted secondary
224 structural characteristics very similar to their homologs in *Saccharomyces cerevisiae*
225 (Figure 4a). However, the structure of the *TpAlg7*-encoding locus was altered as a result
226 of the re-annotation effort and *TpAlg14* is the product of a newly identified gene, which

227 might have prevented the original identification of the pathway. Therefore, *Theileria*
228 parasites likely have a minimal N-glycosylation system. Interestingly, we can find Alg14
229 orthologs by blastp search in *T. orientalis* (TOT_010000184), *T. equi* (BEWA_032670),
230 but not in *T. annulata*. Using the adjacent gene, EngB, as a marker, a look at the *T.*
231 *annulata* genomic region that is syntenic to *TpAlg14* revealed that *T. annulata* has a
232 hypothetical gene annotated on the opposite strand (Figure 4b), which could be an
233 incorrect annotation. A tblastn search of the *T. annulata* genome using *TpAlg14* led to the
234 discovery of a nucleotide sequence which translated results in an alignment with E-value
235 of 7×10^{-15} and 70% identity over the length of the protein, suggesting the existence of an
236 *T. annulata* Alg14 ortholog (*TaAlg14*). In fact, the gene model that was at the *TpAlg14*
237 locus in the original annotation, TP01_0196, was likely a result of an incorrect annotation
238 transfer from *T. annulata* (or vice-versa), since TP01_0196 shared 52% identity with the
239 gene annotated on the opposite strand at the putative *TaAlg14* locus (E-value 4×10^{-131}).
240 Since previous studies have used *T. annulata* as a model *Theileria* parasite, this could be
241 the reason that N-glycosylation was not discovered in this parasite genus.

242 While the presence of N-glycans in *Plasmodium* parasite proteins was initially
243 controversial [34], more recent work provided evidence of short N-glycans on the
244 exterior of *P. falciparum* schizonts and trophozoites [35]. As a key difference,
245 *Plasmodium* parasites have a clear ortholog of the oligosaccharyl transferase STT3 (EC
246 2.4.99.18, PF3D7_1116600 in *P. falciparum* 3D7), which catalyzes the transfer of GlcNAc and
247 GlcNAc₂ to asparagine residues in nascent proteins, and recent work has identified several
248 other proteins in this protein complex in *Plasmodium* genomes [34]. No such ortholog
249 was found in *T. parva* Muguga or *T. annulata* Ankara by blastp or tblastn searches with
250 the *Plasmodium* protein. Since there are STT3 orthologs in *T. equi* and *T. haneyi* (Figure

251 4c), as well as *Cytauxoon felis*, it appears that the absence of STT3 in *T. parva* and *T.*
252 *annulata* represents evolutionary loss of STT3 orthologs in this lineage. This means that
253 while lipid precursor N-glycosylation does likely occur at the ER in these two species,
254 the canonical mechanism of N-glycan precursor transfer to proteins is apparently absent.

255 **Discussion**

256 The re-annotation of the *T. parva* genome has resulted in significant improvement
257 to the accuracy of gene models, showing that this genome is even more gene-dense than
258 previously thought, with the addition of 2.4% of the genome to CDSs as well as the
259 discovery of additional overlapping genes. Multi-gene families appear to have played a
260 prominent role in the evolution of the lineage leading to *T. parva* and *T. annulata* [36],
261 implying a role for these genes in host-pathogen interactions. These genes have
262 diversified and/or expanded in copy number, possibly as an adaptation to a particular
263 niche, since the high density of the genome is strongly suggestive of selection against
264 non-functional DNA. We now have a clearer picture of the structure, copy number, and
265 relative expression level of these genes. In addition, a recently generated sporozoite and
266 sporoblast datasets opens up new opportunities to study differential gene expression
267 throughout other stages of the parasite life cycle [27].

268 The model of transcription that emerges from these recent studies is one of
269 ubiquitous transcription of most genes in the schizont stage, but with a wide range of
270 expression levels [18, 26, 27], suggesting that there are likely important *cis* regulatory
271 motifs that control the level of expression or mRNA stability [18, 37]. Transcription can
272 also arise from potential bidirectional and cryptic promoters with highly prevalent
273 antisense transcription. It remains to be determined if sense and anti-sense transcripts are

274 generated in the same or different cells in culture, an issue that may be addressed with
275 single-cell RNAseq. Due to the short-read nature of our sequencing platform, we were
276 only able to accurately annotate the most prevalent isoform of each gene. The sequencing
277 of full-length transcripts, for example with Pacific Biosciences sequencing technology,
278 would provide a more comprehensive description of the *T. parva* transcriptome,
279 including alternatively spliced variants and the boundaries of overlapping transcripts.

280 In yeast and humans, antisense transcription, defined by the existence of non-
281 coding RNA encoded on the DNA strand opposite to, and overlapping with, that
282 encoding the mRNA, is rare compared to sense transcription [38]. In *T. parva*, however,
283 antisense transcription is highly prevalent throughout the genome (Figure 3b), as has
284 been found in *P. falciparum*, where antisense transcription is synthesized largely by RNA
285 polymerase II [39, 40] and can alter the expression of multigene family members by
286 regulating the packaging of these loci into chromatin [41]. Most of the antisense
287 transcripts seem to completely overlap with their sense counterparts, although the
288 functional relevance of this observation has yet to be determined.

289 The discovery of evidence that N-glycosylation may occur in *Theileria* parasites
290 could open up novel treatment options against *Theileria* infections. N-glycosylation is
291 thought to be important for *Toxoplasma gondii* invasion, growth, and motility [42-44].
292 While the results are somewhat confounded by a lack of inhibitor specificity, treatment
293 with the N-glycosylation inhibitor tunicamycin results in parasites with abnormal
294 endoplasmic reticulum, malformed nuclei, and impaired secretory organelles [45]. While
295 once controversial due to differences in analytical methods, parasite life-cycle stages, and
296 host contamination, *P. falciparum* is now thought to have N-glycosylated proteins,

297 although this is not as frequent a mechanism of protein modification as
298 glycosylphosphatidylinositol [46] . This work has been supported by bioinformatic
299 analyses, finding that *P. falciparum* contains (albeit few) glycosyltransferases [47]. Early
300 work using N-glycosylation inhibitors has shown strong *in vitro* growth inhibition of
301 *Plasmodium* asexual blood stages [48-52], but the function of N-glycosylation of
302 apicomplexan parasite proteins is a topic that requires further study. Importantly, the lack
303 of an STT3 ortholog in *T. parva*, if true, would suggest that protein-targeted N-
304 glycosylation does not occur in this parasite (as does in *Plasmodium*), and may only
305 occur on the ER and potentially the surface of the parasite. Even though cytoplasmic N-
306 glycosyltransferases have been found in bacteria, they have not been found in eukaryotes,
307 and their presence in *T. parva* seems unlikely. The absence of a N-glycan protein transfer
308 system is largely supported by genome-wide searches for the enrichment of N-glycan
309 acceptor sites in *T. annulata* [53]. While N-glycosylation is often touted as an ‘essential’
310 protein modification in eukaryotes [54], the absence of an STT3 ortholog in some
311 *Theileria* species suggests that this process may be critical as a lipid, rather than protein,
312 modification. This does not diminish the potential relevance of N-glycosylation in these
313 parasites. Regardless of whether these short N-glycans provoke host immune responses
314 or play a homeostatic role in parasite protein folding, they could be important therapeutic
315 targets. Finally, given the possibility that glycans encode immunological ‘self’, ‘non-self’
316 or ‘damage’ identities [55], it is tempting to speculate that the absence of proteinaceous
317 N-glycans in *Theileria* species could represent an evolutionary adaptation to immune
318 evasion in a parasite lineage that resides free in the host cytoplasmic environment.
319

320 **Conclusions**

321 This study emphasizes the critical interplay between genome annotations and our
322 knowledge of pathogen biology. The significant improvement of the *T. parva* Muguga
323 reference genome gene annotation will facilitate numerous studies of this parasite to
324 come, and has already given better resolution to genome-wide patterns of gene
325 transcription, including antisense transcription and transcription from multi-gene
326 families. The better the resolution at which we understand gene structure and expression,
327 the more accurately we can characterize and study gene function, novel druggable
328 pathways suitable for interventions and, ultimately, the biology of the pathogen in its
329 different host organisms. For example, the discovery of N-glycosylation precursors in
330 some *Theileria* parasites in the absence of a protein transfer system opens up new
331 questions about the role of lipid N-glycosylation precursors in eukaryote biology as well
332 as the potential evolutionary reasons why protein N-glycosylation would be lost in this
333 apicomplexan lineage.

334

335 **Methods**

336 1) RNA sequencing and genome annotation

337 An RNA sample was obtained from the reference *T. parva* isolate (Muguga) from
338 the haploid schizont stage of the parasite life cycle, which proliferates in host
339 lymphocytes. The extraction method included complement lysis of schizont-infected host
340 lymphocytes, DNase digestion of contaminating host DNA and differential centrifugation
341 to enrich for schizonts [26, 56]. PolyA-enriched RNA was sequenced using Illumina

342 sequencing technology, to produce strand-specific RNAseq data. RNAseq reads were
343 aligned with TopHat and RPKM values calculated using an in-house Perl script.

344

345 2) Genome re-annotation

346 For the re-annotation of the *Theileria parva* genome, a number of evidence tracks
347 were generated and loaded into the genome browser JBrowse [57] for manual curation
348 using the WebApollo plugin [58]. RNAseq reads were aligned to the genome with
349 TopHat [59], a splice-aware alignment tool (Supplementary Table S1, Additional File 1).
350 These alignments were used to generate strand-specific read alignment coverage glyphs
351 and XY plots for visualization in WebApollo. TopHat alignment also yields a file of all
352 reported splice junctions using segmented mapping and coverage information, which is
353 useful for curating intron splice sites. RNAseq reads were also assembled into transcripts
354 using CuffLinks [60] and mapped to the genome with TopHat. We also generated two
355 genome-dependent Trinity/PASA [61] transcriptome assemblies (one reference
356 annotation-dependent and one independent of the reference annotation), as well as one
357 completely *de novo* Trinity transcriptome assembly. A variety of other proteome data
358 were aligned to the genome with AAT [62] and used as evidence tracks, including
359 previously generated *Theileria annulata* mass spectrometry data [63], and all non-
360 *Theileria* apicomplexan proteins from NCBI's RefSeq.

361 In order to assess gene prediction accuracy before the manual curation phase, a set
362 of 342 high-confidence *T. parva* gene models were selected from the current reference
363 annotation on the basis of two criteria: (1) RNAseq reads must cover each exon in the
364 gene, (2) Trinity *de novo* assembled transcripts and read coverage must be concordant

365 with the presence or absence of any introns in the gene model. Out of these 342 genes, 50
366 were randomly selected as a validation set and the remaining 292 were used as a training
367 gene set for gene prediction software. The exon distribution of the validation set closely
368 resembles that of the training set (Supplementary Table S8, Additional File 1).

369 Multiple gene prediction software tools were used and then assessed by the
370 accuracy with which they predict the validation set using an in-house script. These
371 included: *i*) Augustus [64], using RNAseq reads, the *T. parva* training gene set, or no
372 evidence; *ii*) Semi-HMM-based Nucleic Acid Parser (SNAP) [65] and Glimmer [66]
373 were trained with the *T. parva* training set; *iii*) Fgenesh [67] used a pre-existing training
374 set of *Plasmodium* genes from its website; *iv*) the *ab initio* predictor GeneMark-ES [68].
375 Finally, gene models were selected with the consensus predictor Evidence Modeler
376 (EVM) [69], using 57, differently-weighted combinations of the other evidence, while
377 maximizing prediction accuracy (Supplementary Figure S8, Additional File 1). Based on
378 their performance in comparison with the validation set, only the top four EVM
379 predictions were loaded as evidence tracks for use in manual curation (Supplementary
380 Figure S9, Additional File 1). tRNA and rRNA predictions were generated using
381 tRNAscan-SE [70] and RNAmmer [71] and loaded as evidence tracks, along with the
382 original *T. parva* Muguga annotation (Supplementary Figure S10, Additional File 1). A
383 genome-wide, double-pass, manual curation of all gene models was completed, weighing
384 the RNAseq evidence over the evidence from alignments with homologs from other
385 species and the gene prediction programs. The annotation assignments were allocated in
386 50 kb segments, with different annotators doing adjacent segments, as well as altering the
387 annotator for the first and second pass in order to reduce annotator bias.

388 Functional annotation of the *T. parva* proteome consisted of HMM3 searches of
389 the complete proteome against our custom HMM collection that includes TIGRFams
390 [72], Pfams [73], as well as custom-built HMMs [74] and RAPSearch2 searches against
391 UniRef100 (with a cutoff of 1×10^{-10}). In addition, a TMHMM search which was used to
392 assign "putative integral membrane protein" to proteins with 3 or more helical spans
393 (assuming there were no other hits to the previous searches). These searches were
394 synthesized using Attributor (<https://github.com/jorvis/Attributor>) to generate the final
395 annotation based on the different evidence sources to assign gene product names, EC
396 numbers, GO terms and gene symbols to genes, conservatively where possible.

397

398 3) Multi-gene family clustering

399 Genes were clustered with OrthoMCL, using an inflation value of 4 and a BLAST
400 p-value cutoff of 10^{-5} , as previously done [19]. All individual conserved domain searches
401 were done using NCBI's Conserved Domain Database version 3.11 [75] with 45,746
402 PSSMs, with an E-value threshold of 0.01 and a composition based statistics adjustment.
403 HMM searches of the entire PFAM database were done using default settings.

404

405 **List of abbreviations**

406 ECF: East Coast fever

407 UTR: Untranslated Regions

408 CDS: Coding Sequence

409 RPKM: reads per kilobase million

410 Tpr: *T. parva* repeat family

411 SVSP: subtelomeric variable secreted protein

412 SNAP: Semi-HMM-based Nucleic Acid Parser

413

414 **Ethics approval and consent to participate**

415 Not applicable.

416

417 **Consent for publication**

418 Not applicable.

419

420 **Availability of data and materials**

421 The *T. parva* Muguga re-annotation is publicly available in GenBank and can be

422 visualized at the following online link (http://jbrowse.igs.umaryland.edu/t_parva/), as

423 well as under the NCBI BioProject PRJNA16138.

424

425 **Competing interests**

426 The authors declare that they have no competing interests.

427

428 **Funding**

429 This work was supported in part by an Immunity and Infection T32 training grant,

430 NIH/NIAID T32 AI007540-14, by the by the Bill and Melinda Gates Foundation

431 (OPP1078791), and by the United States Department of Agriculture, Agricultural

432 Research Service (USDA-ARS) through agreement #59-5348-4-001.

433

434 **Authors' contributions**

435 RP isolated *T. parva* schizont RNA for RNAseq using differential centrifugation and
436 standard kits as described. JO, OOI, and PK built and maintained the JBrowse instance
437 for manual curation using the WebApollo plugin, as well as the functional annotation
438 pipeline. KT, JO, HTG, OOI, PK, and SBAI generated alignment tracks to assist the
439 annotation. KT, HTG, OOI, NCP, SBAI, JCS completed manual re-annotation of the
440 genome. KT performed all other analyses. JCS, RPB, CD, VMN, and LF conceived the
441 study design. KT and JCS wrote the manuscript. All authors critically reviewed and
442 approved the manuscript.

443

444 **Acknowledgements**

445 We would like to thank Donald P. Knowles for his kind and helpful feedback and support
446 for this project.

447

448

449

450

451

452

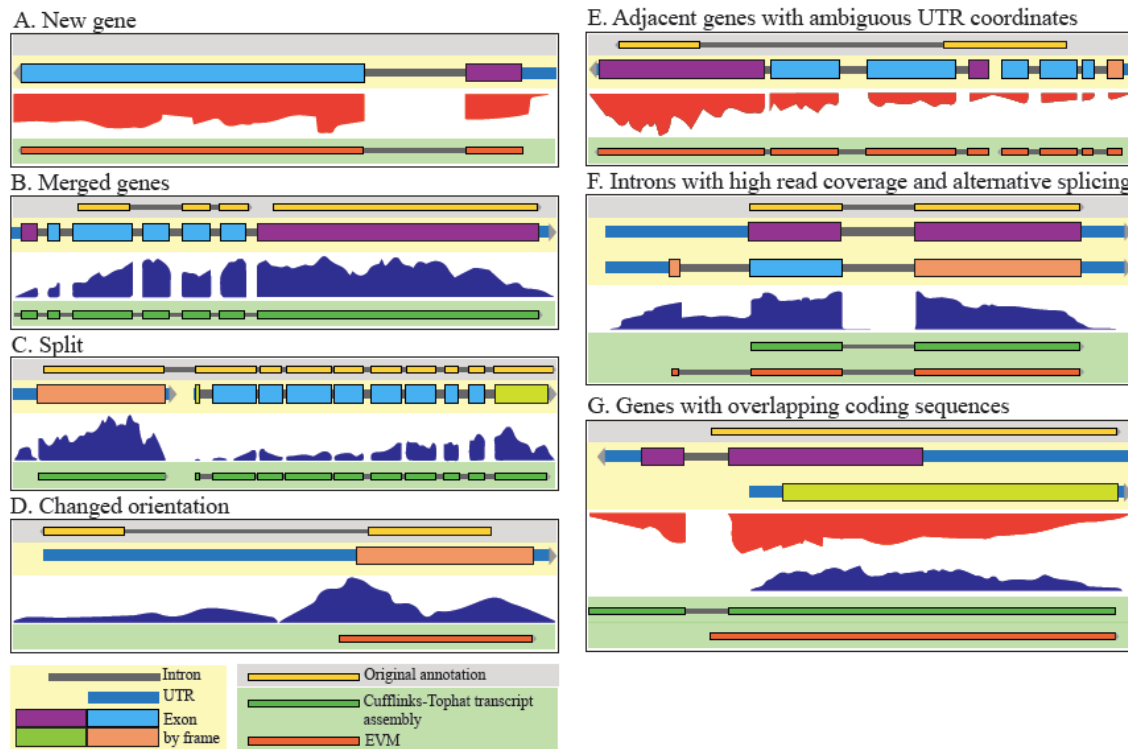
453

454

455

456

457 **Figures and Figure Legends**



458

459 **Figure 1. Manual gene model curation examples.** Several tracks are shown: updated
460 gene model (beige background), original (2005) gene annotation (grey background),
461 RNAseq data (white background), transcript assembly (dark green, on green
462 background), and EVM predictions (orange, on green background). (A) A new gene
463 discovered on the basis of RNAseq data (TpMuguga_03g02005). (B) A case where two
464 genes in the 2005 annotation merge in the new annotation on the basis of RNAseq read
465 coverage (TpMuguga_04g02435). (C) A case where a gene in the 2005 annotation has
466 been split into two genes in the new annotation (TpMuguga_04g02190 and
467 TpMuguga_04g02185). (D) A case where a gene has been reversed in orientation on the
468 basis of RNAseq data (TpMuguga_02g02095). (E) A case where overlapping genes led
469 to ambiguity in UTR coordinates, and so the UTRs were not defined in this intergenic
470 region (TpMuguga_01g00527 and TpMuguga_01g00528). (F) A case of a single gene

471 where alternative splicing exists (as seen by significant read coverage in at least one
472 intronic region), but there is one most prevalent isoform (TpMuguga_03g00622). (G) A
473 case of two genes that overlap by coding sequences. Coding exons are colored by reading
474 frame (TpMuguga_05g00017 and TpMuguga_05g00018).

475

476

477

478

479

480

481

482

483

484

485

486

487

488

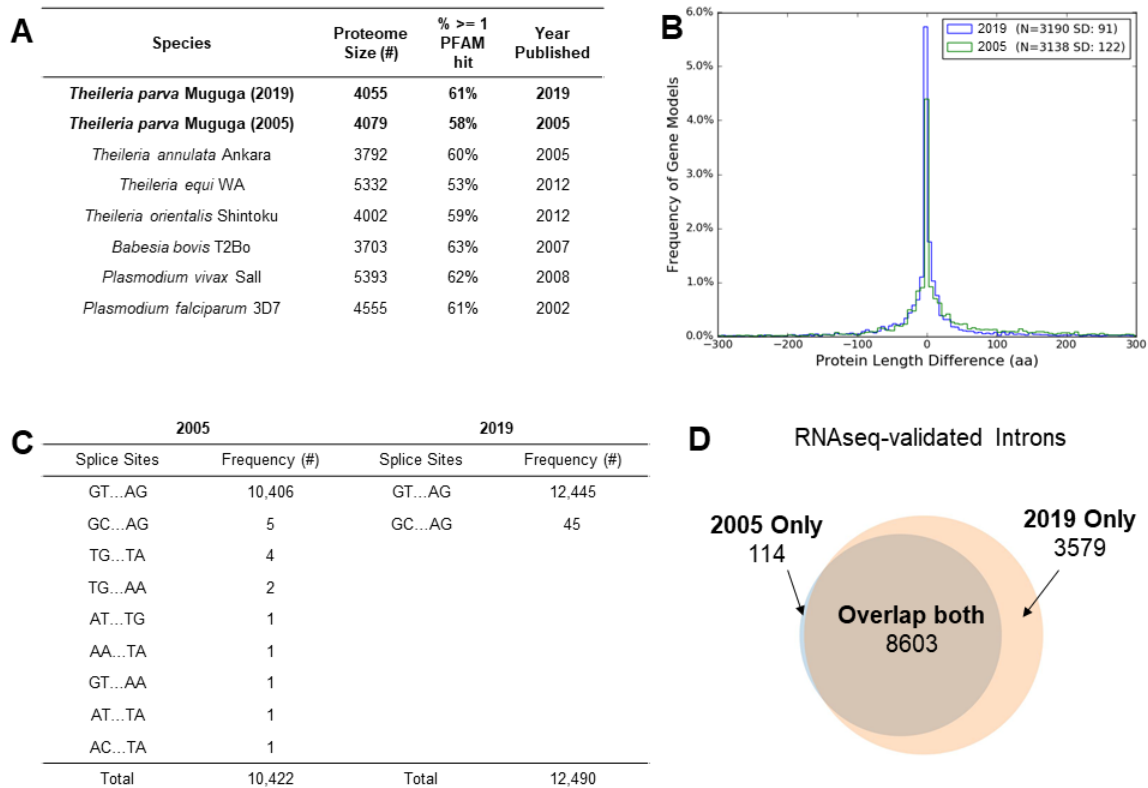
489

490

491

492

493



494

495 **Figure 2. Comparative metrics of original and new *T. parva* annotations.**

496 (A) The percentage of proteins with at least one PFAM domain found by Hidden Markov
 497 Model searches of the predicted proteomes of the new *T. parva* Muguga annotation was
 498 2% higher than those in the 2005 annotation, implying that the new annotation captures
 499 functional elements that were previously missed. (B) The new *T. parva* Muguga
 500 annotation has more reciprocal best-hit orthologs (N) with *T. annulata* Ankara than the
 501 2005 *T. parva* Muguga annotation. The variation in protein length (SD) between *T. parva*
 502 and *T. annulata* ortholog pairs is greatly reduced in the new relative to the original *T.*
 503 *parva* annotation. Only nuclear genes were used for this analysis. The x-axis was limited
 504 to the range -300 to +300 for easy visual interpretation. (C) The number of canonical
 505 GT/AG intron splice sites increased and the number of non-canonical intron splice site

506 combinations decreased in the new *T. parva* Muguga annotation compared to the 2005
507 annotation. (D) The number and proportion of introns validated by at least one RNAseq
508 read increased in the new *T. parva* Muguga annotation compared to the 2005 annotation.
509 These lines of evidence suggest that the new annotation is more accurate, and also
510 considerably more consistent with the RNAseq data, as expected.

511

512

513

514

515

516

517

518

519

520

521

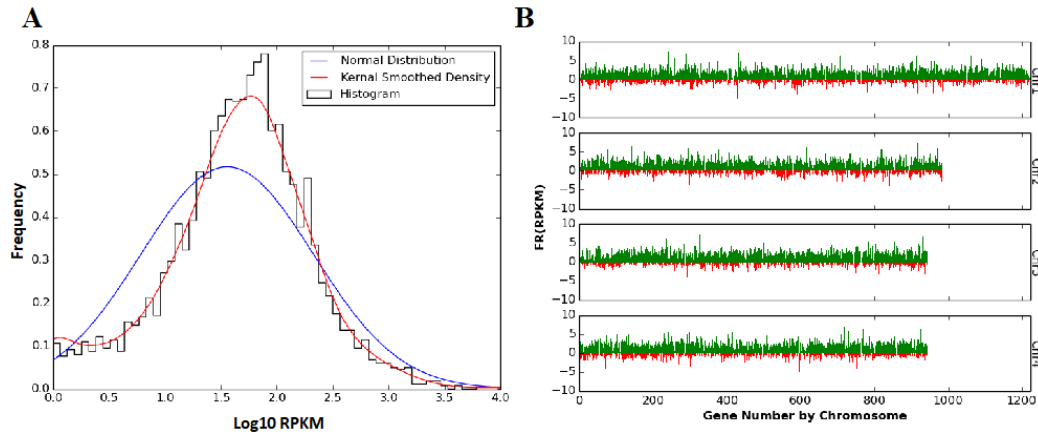
522

523

524

525

526



527

528 **Figure 3. Distribution of RNAseq RPKM values for *T. parva* Muguga genes (A) A**

529 histogram of sense RPKM values after logarithmic transformation of the data.

530 Frequencies on the y-axis correspond to probability density. The blue line shows a

531 normal distribution around the same median, while the red line shows a more reliable

532 fixed-width, Gaussian, kernel-smoothed estimate of the probability density. (B) The

533 sense (green) and antisense (red) reads per kilobase transcript per million reads (RPKM)

534 after fourth-root transformation of the data. Genes are sorted by position on the

535 chromosome for all four nuclear chromosomes of *T. parva* Muguga.

536

537

538

539

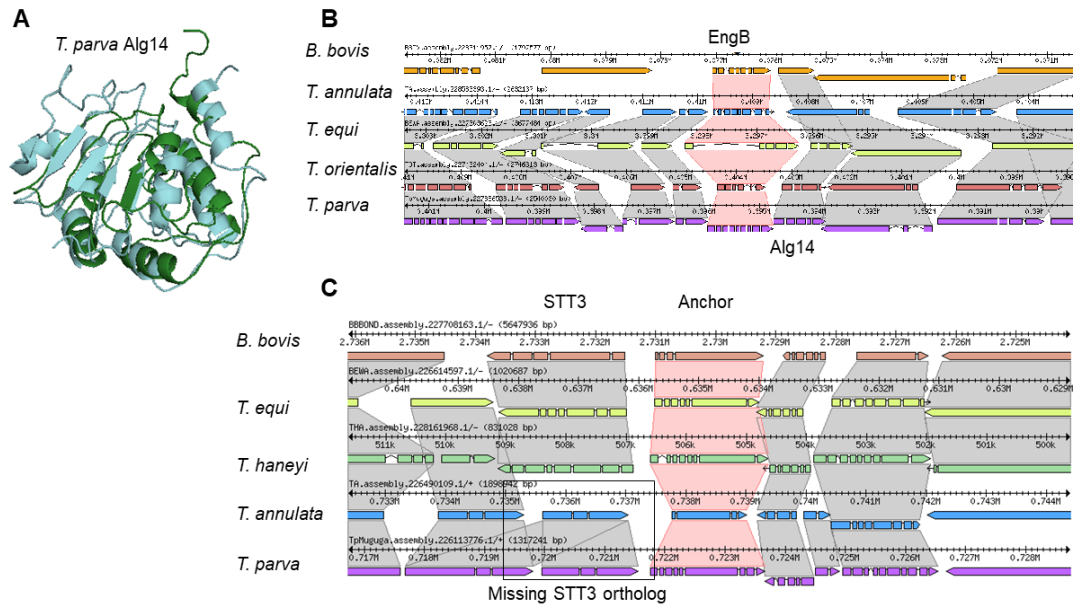
540

541

542

543

544



545

546 **Figure 4. The uncovered *Theileria parva* Alg14 shows a similar predicted structure**

547 **to the empirically determined *Saccharomyces cerevisiae* Alg14 protein structure, and**

548 **is syntenic in multiple piroplasm. (A) A Phyre2 prediction of *T. parva* Alg14**

549 **(TpAlg14; green; TpMuguga_01g02045) and the Protein Database (<http://www.rcsb.org/>)**

550 **[76] nuclear magnetic resonance structure of *Saccharomyces cerevisiae* Alg14 (ScAlg14;**

551 **teal; PDB 2JZC) were aligned in MacPyMol (<https://pymol.org/2/>) [77]. (B) Shown are**

552 **the syntenic regions around Alg14 orthologs (synteny in grey), using the adjacent gene**

553 **EngB as an anchor (synteny in red) in the Sybil software package [78]. (C) Shown are the**

554 **syntenic regions around STT3 orthologs (synteny in grey), using a *B. bovis* STT3-**

555 **adjacent gene ([BBOV_I1000210](http://www.ncbi.nlm.nih.gov/nuccore/BBOV_I1000210)) as an anchor (synteny in red) in the Sybil software**

556 **package.**

557

558

559

560

561 **References**

- 562 1. Spielman DJ: **XVI Public-Private Partnerships and Pro-Poor Livestock**
563 **Research: The Search for an East Coast Fever Vaccine**, vol. 1. Washington,
564 D.C.: The National Academies Press; 2009.
- 565 2. Herrero M, Thornton PK, Notenbaert AM, Wood S, Msangi S, Freeman HA,
566 Bossio D, Dixon J, Peters M, van de Steeg J *et al*: **Smart investments in**
567 **sustainable food production: revisiting mixed crop-livestock systems**. *Science*
568 2010, **327**(5967):822-825.
- 569 3. Nkedianye D, Radeny M, Kristjanson P, Herrero M: **Assessing returns to land**
570 **and changing livelihood strategies in Kitengela**. In: *Staying Maasai?*
571 *Livelihoods, conservation and development in East African Rangelands*. Edited
572 by Homewood K, Kristjanson P, Chevenix Trench P. Dordrecht, The Netherlands:
573 Springer; 2009: 115-150.
- 574 4. Baldwin CL, Black SJ, Brown WC, Conrad PA, Goddeeris BM, Kinuthia SW,
575 Lalor PA, MacHugh ND, Morrison WI, Morzaria SP *et al*: **Bovine T cells, B**
576 **cells, and null cells are transformed by the protozoan parasite Theileria**
577 **parva**. *Infection and immunity* 1988, **56**(2):462-467.
- 578 5. Tindih HS, Geysen D, Goddeeris BM, Awino E, Dobbelaere DA, Naessens J: **A**
579 **Theileria parva isolate of low virulence infects a subpopulation of**
580 **lymphocytes**. *Infection and immunity* 2012, **80**(3):1267-1273.
- 581 6. Irvin AD, Mwamachi DM: **Clinical and diagnostic features of East Coast fever**
582 **(Theileria parva) infection of cattle**. *The Veterinary record* 1983, **113**(9):192-
583 198.

- 584 7. Fry LM, Schneider DA, Frevert CW, Nelson DD, Morrison WI, Knowles DP:
585 **East Coast Fever Caused by *Theileria parva* Is Characterized by**
586 **Macrophage Activation Associated with Vasculitis and Respiratory Failure.**
587 *PLoS One* 2016, **11**(5):e0156004.
- 588 8. Hotez PJ, Molyneux DH, Fenwick A, Kumaresan J, Sachs SE, Sachs JD, Savioli
589 **L: Control of neglected tropical diseases.** *The New England journal of medicine*
590 2007, **357**(10):1018-1027.
- 591 9. Reed SL, McKerrow JH: **Why Funding for Neglected Tropical Diseases**
592 **Should Be a Global Priority.** *Clin Infect Dis* 2018, **67**(3):323-326.
- 593 10. Sette A, Rappuoli R: **Reverse vaccinology: developing vaccines in the era of**
594 **genomics.** *Immunity* 2010, **33**(4):530-541.
- 595 11. Seib KL, Dougan G, Rappuoli R: **The key role of genomics in modern vaccine**
596 **and drug design for emerging infectious diseases.** *PLoS genetics* 2009,
597 **5**(10):e1000612.
- 598 12. Hotez PJ, Fenwick A, Ray SE, Hay SI, Molyneux DH: **"Rapid impact" 10 years**
599 **after: The first "decade" (2006-2016) of integrated neglected tropical disease**
600 **control.** *PLoS Negl Trop Dis* 2018, **12**(5):e0006137.
- 601 13. Chaudhary K, Roos DS: **Protozoan genomics for drug discovery.** *Nature*
602 *biotechnology* 2005, **23**(9):1089-1091.
- 603 14. Oberg AL, Kennedy RB, Li P, Ovsyannikova IG, Poland GA: **Systems biology**
604 **approaches to new vaccine development.** *Current opinion in immunology* 2011,
605 **23**(3):436-443.

- 606 15. Wakaguri H, Suzuki Y, Sasaki M, Sugano S, Watanabe J: **Inconsistencies of**
607 **genome annotations in apicomplexan parasites revealed by 5'-end-one-pass**
608 **and full-length sequences of oligo-capped cDNAs.** *BMC Genomics* 2009,
609 **10**:312.
- 610 16. Yeoh LM, Lee VV, McFadden GI, Ralph SA: **Alternative Splicing in**
611 **Apicomplexan Parasites.** *MBio* 2019, **10**(1).
- 612 17. Yandell M, Ence D: **A beginner's guide to eukaryotic genome annotation.**
613 *Nature reviews Genetics* 2012, **13**(5):329-342.
- 614 18. Tretina K, Pelle R, Silva JC: **Cis regulatory motifs and antisense**
615 **transcriptional control in the apicomplexan *Theileria parva*.** *BMC genomics*
616 2016, **17**(1):128.
- 617 19. Gardner MJ, Bishop R, Shah T, de Villiers EP, Carlton JM, Hall N, Ren Q,
618 Paulsen IT, Pain A, Berriman M *et al*: **Genome sequence of *Theileria parva*, a**
619 **bovine pathogen that transforms lymphocytes.** *Science* 2005, **309**(5731):134-
620 137.
- 621 20. Pain A, Renauld H, Berriman M, Murphy L, Yeats CA, Weir W, Kerhornou A,
622 Aslett M, Bishop R, Bouchier C *et al*: **Genome of the host-cell transforming**
623 **parasite *Theileria annulata* compared with *T. parva*.** *Science* 2005,
624 **309**(5731):131-133.
- 625 21. Collins L, Penny D: **Complex spliceosomal organization ancestral to extant**
626 **eukaryotes.** *Molecular biology and evolution* 2005, **22**(4):1053-1066.

- 627 22. Russell K, Hasenkamp S, Emes R, Horrocks P: **Analysis of the spatial and**
628 **temporal arrangement of transcripts over intergenic regions in the human**
629 **malarial parasite Plasmodium falciparum.** *BMC genomics* 2013, **14**:267.
- 630 23. Unlu AH, Tajeri S, Bilgic HB, Eren H, Karagenc T, Langsley G: **The secreted**
631 **Theileria annulata Ta9 protein contributes to activation of the AP-1**
632 **transcription factor.** *PLoS One* 2018, **13**(5):e0196875.
- 633 24. Painter HJ, Carrasquilla M, Llinas M: **Capturing in vivo RNA transcriptional**
634 **dynamics from the malaria parasite Plasmodium falciparum.** *Genome*
635 *research* 2017.
- 636 25. Rovira-Graells N, Gupta AP, Planet E, Crowley VM, Mok S, Ribas de Pouplana
637 L, Preiser PR, Bozdech Z, Cortes A: **Transcriptional variation in the malaria**
638 **parasite Plasmodium falciparum.** *Genome research* 2012, **22**(5):925-938.
- 639 26. Bishop R, Shah T, Pelle R, Hoyle D, Pearson T, Haines L, Brass A, Hulme H,
640 Graham SP, Taracha EL *et al*: **Analysis of the transcriptome of the protozoan**
641 **Theileria parva using MPSS reveals that the majority of genes are**
642 **transcriptionally active in the schizont stage.** *Nucleic acids research* 2005,
643 **33**(17):5503-5511.
- 644 27. Tonui T, Corredor-Moreno P, Kanduma E, Njuguna J, Njahira MN, Nyanjom SG,
645 Silva JC, Djikeng A, Pelle R: **Transcriptomics reveal potential vaccine**
646 **antigens and a drastic increase of upregulated genes during Theileria parva**
647 **development from arthropod to bovine infective stages.** *PLoS One* 2018,
648 **13**(10):e0204047.

- 649 28. Deitsch KW, Dzikowski R: **Variant Gene Expression and Antigenic Variation**
650 **by Malaria Parasites.** *Annu Rev Microbiol* 2017, **71**:625-641.
- 651 29. Bishop R, Musoke A, Morzaria S, Sohanpal B, Gobright E: **Concerted evolution**
652 **at a multicopy locus in the protozoan parasite *Theileria parva*: extreme**
653 **divergence of potential protein-coding sequences.** *Mol Cell Biol* 1997,
654 **17(3):1666-1673.**
- 655 30. Schmuckli-Maurer J, Casanova C, Schmied S, Affentranger S, Parvanova I,
656 Kang'a S, Nene V, Katzer F, McKeever D, Muller J *et al*: **Expression analysis of**
657 **the *Theileria parva* subtelomere-encoded variable secreted protein gene**
658 **family.** *PLoS One* 2009, **4(3):e4839.**
- 659 31. Marsolier J, Perichon M, DeBarry JD, Villoutreix BO, Chluba J, Lopez T,
660 Garrido C, Zhou XZ, Lu KP, Fritsch L *et al*: ***Theileria* parasites secrete a prolyl**
661 **isomerase to maintain host leukocyte transformation.** *Nature* 2015,
662 **16(520):378-382.**
- 663 32. Schwarz F, Aebi M: **Mechanisms and principles of N-linked protein**
664 **glycosylation.** *Curr Opin Struct Biol* 2011, **21(5):576-582.**
- 665 33. Samuelson J, Robbins PW: **Effects of N-glycan precursor length diversity on**
666 **quality control of protein folding and on protein glycosylation.** *Semin Cell*
667 *Dev Biol* 2015, **41**:121-128.
- 668 34. Tamana S, Promponas VJ: **An updated view of the oligosaccharyltransferase**
669 **complex in *Plasmodium*.** *Glycobiology* 2019, **29(5):385-396.**
- 670 35. Bushkin GG, Ratner DM, Cui J, Banerjee S, Duraisingh MT, Jennings CV,
671 Dvorin JD, Gubbels MJ, Robertson SD, Steffen M *et al*: **Suggestive evidence for**

- 672 **Darwinian Selection against asparagine-linked glycans of Plasmodium**
673 **falciparum and Toxoplasma gondii.** *Eukaryot Cell* 2010, **9**(2):228-241.
- 674 36. Reid AJ: **Large, rapidly evolving gene families are at the forefront of host-**
675 **parasite interactions in Apicomplexa.** *Parasitology* 2015, **142 Suppl 1**:S57-70.
- 676 37. Pieszko M, Weir W, Goodhead I, Kinnaird J, Shiels B: **ApiAP2 Factors as**
677 **Candidate Regulators of Stochastic Commitment to Merozoite Production in**
678 **Theileria annulata.** *PLoS neglected tropical diseases* 2015, **9**(8):e0003933.
- 679 38. Ozsolak F, Kapranov P, Foissac S, Kim SW, Fishilevich E, Monaghan AP, John
680 B, Milos PM: **Comprehensive polyadenylation site maps in yeast and human**
681 **reveal pervasive alternative polyadenylation.** *Cell* 2010, **143**(6):1018-1029.
- 682 39. Militello KT, Patel V, Chessler AD, Fisher JK, Kasper JM, Gunasekera A, Wirth
683 DF: **RNA polymerase II synthesizes antisense RNA in Plasmodium**
684 **falciparum.** *RNA* 2005, **11**(4):365-370.
- 685 40. Lopez-Barragan MJ, Lemieux J, Quinones M, Williamson KC, Molina-Cruz A,
686 Cui K, Barillas-Mury C, Zhao K, Su XZ: **Directional gene expression and**
687 **antisense transcripts in sexual and asexual stages of Plasmodium falciparum.**
688 *BMC Genomics* 2011, **12**:587.
- 689 41. Jing Q, Cao L, Zhang L, Cheng X, Gilbert N, Dai X, Sun M, Liang S, Jiang L:
690 **Plasmodium falciparum var Gene Is Activated by Its Antisense Long**
691 **Noncoding RNA.** *Front Microbiol* 2018, **9**:3117.
- 692 42. Fauquenoy S, Morelle W, Hovasse A, Bednarczyk A, Slomianny C, Schaeffer C,
693 Van Dorsselaer A, Tomavo S: **Proteomics and glycomics analyses of N-**

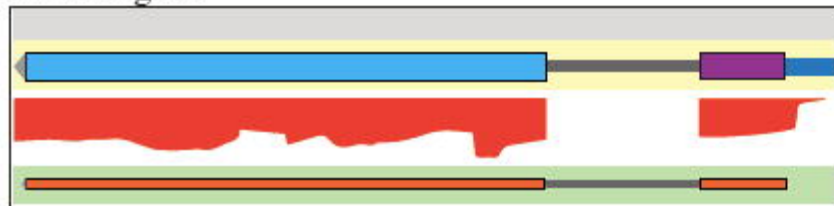
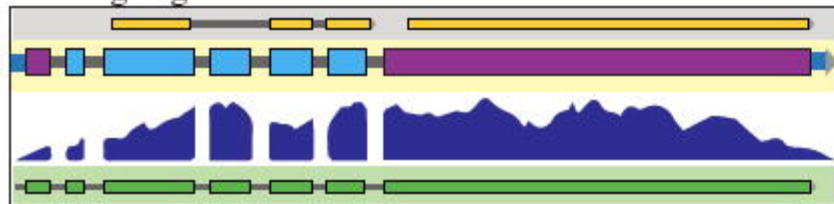
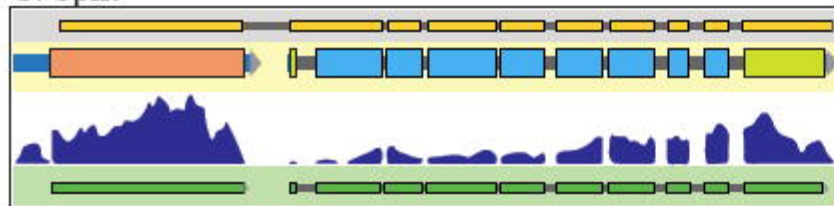
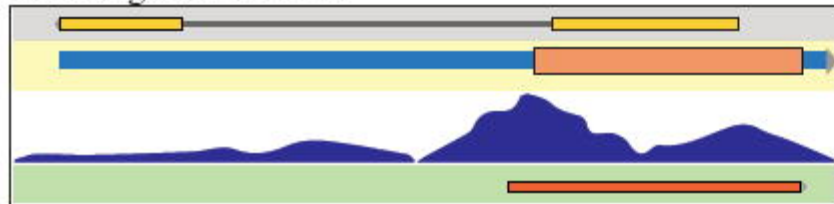
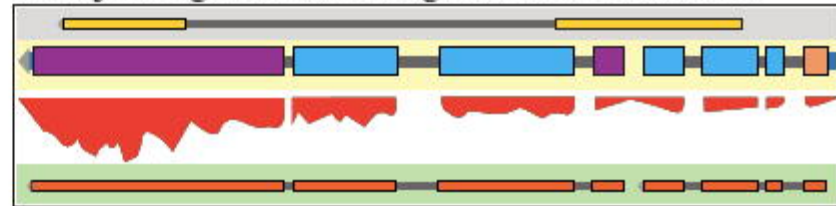
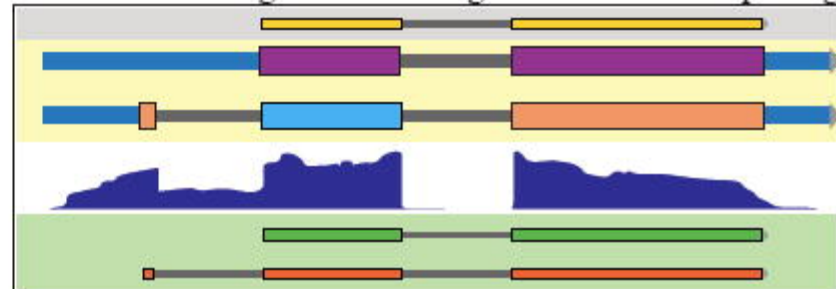
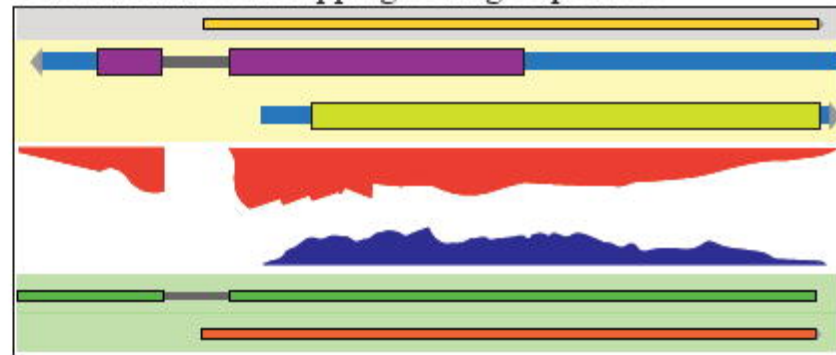
- 694 **glycosylated structures involved in *Toxoplasma gondii*--host cell interactions.**
695 *Mol Cell Proteomics* 2008, **7**(5):891-910.
- 696 43. Gas-Pascual E, Ichikawa HT, Sheikh MO, Serji MI, Deng B, Mandalasi M,
697 Bandini G, Samuelson J, Wells L, West CM: **CRISPR/Cas9 and glycomics tools**
698 **for *Toxoplasma* glycobiology.** *J Biol Chem* 2019, **294**(4):1104-1125.
- 699 44. Sidik SM, Huet D, Ganesan SM, Huynh MH, Wang T, Nasamu AS, Thiru P,
700 Saeij JPJ, Carruthers VB, Niles JC *et al*: **A Genome-wide CRISPR Screen in**
701 ***Toxoplasma* Identifies Essential Apicomplexan Genes.** *Cell* 2016,
702 **166**(6):1423-1435 e1412.
- 703 45. Luk FC, Johnson TM, Beckers CJ: **N-linked glycosylation of proteins in the**
704 **protozoan parasite *Toxoplasma gondii*.** *Molecular and biochemical*
705 *parasitology* 2008, **157**(2):169-178.
- 706 46. Cova M, Rodrigues JA, Smith TK, Izquierdo L: **Sugar activation and**
707 **glycosylation in *Plasmodium*.** *Malar J* 2015, **14**:427.
- 708 47. Samuelson J, Banerjee S, Magnelli P, Cui J, Kelleher DJ, Gilmore R, Robbins
709 PW: **The diversity of dolichol-linked precursors to Asn-linked glycans likely**
710 **results from secondary loss of sets of glycosyltransferases.** *Proceedings of the*
711 *National Academy of Sciences of the United States of America* 2005, **102**(5):1548-
712 1553.
- 713 48. Trigg PI, Hirst SI, Shakespeare PG, Tappenden L: **Labelling of membrane**
714 **glycoprotein in erythrocytes infected with *Plasmodium knowlesi*.** *Bull World*
715 *Health Organ* 1977, **55**(2-3):205-209.

- 716 49. Udeinya IJ, Van Dyke K: **Labelling of membrane glycoproteins of cultivated**
717 **Plasmodium falciparum.** *Bull World Health Organ* 1980, **58**(3):445-448.
- 718 50. Udeinya IJ, Van Dyke K: **Plasmodium falciparum: synthesis of glycoprotein**
719 **by cultured erythrocytic stages.** *Experimental parasitology* 1981, **52**(3):297-
720 302.
- 721 51. Udeinya IJ, Van Dyke K: **Concurrent inhibition by tunicamycin of**
722 **glycosylation and parasitemia in malarial parasites (Plasmodium falciparum)**
723 **cultured in human erythrocytes.** *Pharmacology* 1981, **23**(3):165-170.
- 724 52. Udeinya IJ, Van Dyke K: **2-Deoxyglucose: inhibition of parasitemia and of**
725 **glucosamine incorporation into glycosylated macromolecules, in malarial**
726 **parasites (Plasmodium falciparum).** *Pharmacology* 1981, **23**(3):171-175.
- 727 53. Cui J, Smith T, Robbins PW, Samuelson J: **Darwinian selection for sites of Asn-**
728 **linked glycosylation in phylogenetically disparate eukaryotes and viruses.**
729 *Proc Natl Acad Sci U S A* 2009, **106**(32):13421-13426.
- 730 54. Lombard J: **The multiple evolutionary origins of the eukaryotic N-**
731 **glycosylation pathway.** *Biol Direct* 2016, **11**:36.
- 732 55. Maverakis E, Kim K, Shimoda M, Gershwin ME, Patel F, Wilken R,
733 Raychaudhuri S, Ruhaak LR, Lebrilla CB: **Glycans in the immune system and**
734 **The Altered Glycan Theory of Autoimmunity: a critical review.** *J Autoimmun*
735 2015, **57**:1-13.
- 736 56. Graham SP, Pelle R, Honda Y, Mwangi DM, Tonukari NJ, Yamage M, Glew EJ,
737 de Villiers EP, Shah T, Bishop R *et al*: **Theileria parva candidate vaccine**
738 **antigens recognized by immune bovine cytotoxic T lymphocytes.** *Proceedings*

- 739 *of the National Academy of Sciences of the United States of America* 2006,
740 **103**(9):3286-3291.
- 741 57. Skinner ME, Uzilov AV, Stein LD, Mungall CJ, Holmes IH: **JBrowse: a next-**
742 **generation genome browser.** *Genome research* 2009, **19**(9):1630-1638.
- 743 58. Lee E, Helt GA, Reese JT, Munoz-Torres MC, Childers CP, Buels RM, Stein L,
744 Holmes IH, Elisk CG, Lewis SE: **Web Apollo: a web-based genomic**
745 **annotation editing platform.** *Genome biology* 2013, **14**(8):R93.
- 746 59. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL: **TopHat2:**
747 **accurate alignment of transcriptomes in the presence of insertions, deletions**
748 **and gene fusions.** *Genome biology* 2013, **14**(4):R36.
- 749 60. Roberts A, Pimentel H, Trapnell C, Pachter L: **Identification of novel**
750 **transcripts in annotated genomes using RNA-Seq.** *Bioinformatics* 2011,
751 **27**(17):2325-2329.
- 752 61. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis
753 X, Fan L, Raychowdhury R, Zeng Q *et al*: **Full-length transcriptome assembly**
754 **from RNA-Seq data without a reference genome.** *Nature biotechnology* 2011,
755 **29**(7):644-652.
- 756 62. Huang X, Adams MD, Zhou H, Kerlavage AR: **A tool for analyzing and**
757 **annotating genomic sequences.** *Genomics* 1997, **46**(1):37-45.
- 758 63. Witschi M, Xia D, Sanderson S, Baumgartner M, Wastling JM, Dobbelaere DA:
759 **Proteomic analysis of the *Theileria annulata* schizont.** *International journal for*
760 *parasitology* 2013, **43**(2):173-180.

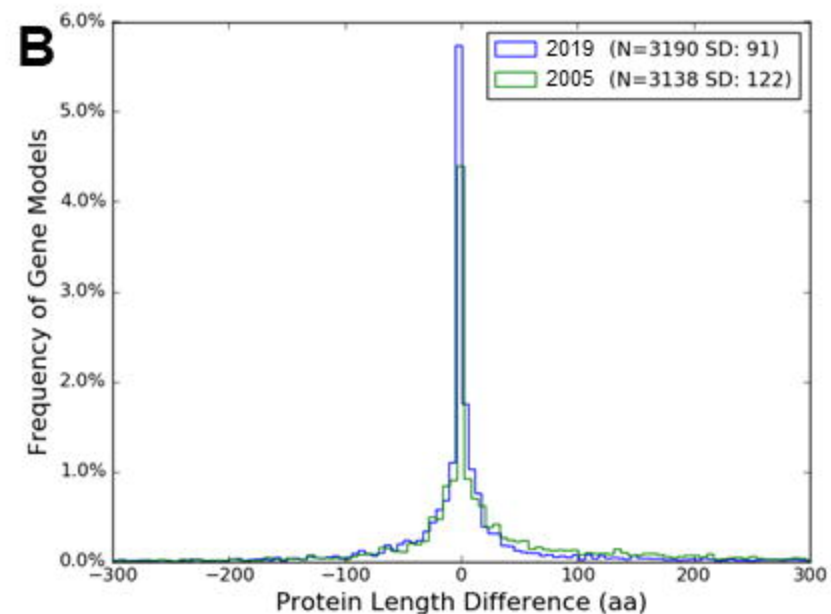
- 761 64. Stanke M, Morgenstern B: **AUGUSTUS: a web server for gene prediction in**
762 **eukaryotes that allows user-defined constraints.** *Nucleic Acids Res* 2005,
763 **33**(Web Server issue):W465-467.
- 764 65. Korf I: **Gene finding in novel genomes.** *BMC bioinformatics* 2004, **5**:59.
- 765 66. Delcher AL, Bratke KA, Powers EC, Salzberg SL: **Identifying bacterial genes**
766 **and endosymbiont DNA with Glimmer.** *Bioinformatics* 2007, **23**(6):673-679.
- 767 67. Solovyev V, Kosarev P, Seledsov I, Vorobyev D: **Automatic annotation of**
768 **eukaryotic genes, pseudogenes and promoters.** *Genome biology* 2006, **7 Suppl**
769 **1**:S10 11-12.
- 770 68. Borodovsky M, Lomsadze A: **Eukaryotic gene prediction using**
771 **GeneMark.hmm-E and GeneMark-ES.** *Current protocols in bioinformatics /*
772 *editorial board, Andreas D Baxevanis [et al]* 2011, **Chapter 4**:Unit 4 6 1-10.
- 773 69. Haas BJ, Salzberg SL, Zhu W, Pertea M, Allen JE, Orvis J, White O, Buell CR,
774 Wortman JR: **Automated eukaryotic gene structure annotation using**
775 **EvidenceModeler and the Program to Assemble Spliced Alignments.** *Genome*
776 *biology* 2008, **9**(1):R7.
- 777 70. Lowe TM, Eddy SR: **tRNAscan-SE: a program for improved detection of**
778 **transfer RNA genes in genomic sequence.** *Nucleic acids research* 1997,
779 **25**(5):955-964.
- 780 71. Lagesen K, Hallin P, Rodland EA, Staerfeldt HH, Rognes T, Ussery DW:
781 **RNAmmer: consistent and rapid annotation of ribosomal RNA genes.** *Nucleic*
782 *acids research* 2007, **35**(9):3100-3108.

- 783 72. Haft DH, Selengut JD, White O: **The TIGRFAMs database of protein families.**
784 *Nucleic acids research* 2003, **31**(1):371-373.
- 785 73. Finn RD, Bateman A, Clements J, Coggill P, Eberhardt RY, Eddy SR, Heger A,
786 Hetherington K, Holm L, Mistry J *et al*: **Pfam: the protein families database.**
787 *Nucleic acids research* 2014, **42**(Database issue):D222-230.
- 788 74. Li L, Stoeckert CJ, Jr., Roos DS: **OrthoMCL: identification of ortholog groups**
789 **for eukaryotic genomes.** *Genome research* 2003, **13**(9):2178-2189.
- 790 75. Marchler-Bauer A, Bo Y, Han L, He J, Lanczycki CJ, Lu S, Chitsaz F, Derbyshire
791 MK, Geer RC, Gonzales NR *et al*: **CDD/SPARCLE: functional classification of**
792 **proteins via subfamily domain architectures.** *Nucleic Acids Res* 2017,
793 **45**(D1):D200-D203.
- 794 76. Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, Shindyalov
795 IN, Bourne PE: **The Protein Data Bank.** *Nucleic Acids Res* 2000, **28**(1):235-242.
- 796 77. Janson G, Zhang C, Prado MG, Paiardini A: **PyMod 2.0: improvements in**
797 **protein sequence-structure analysis and homology modeling within PyMOL.**
798 *Bioinformatics* 2017, **33**(3):444-446.
- 799 78. Riley DR, Angiuoli SV, Crabtree J, Dunning Hotopp JC, Tettelin H: **Using Sybil**
800 **for interactive comparative genomics of microbes on the web.** *Bioinformatics*
801 2012, **28**(2):160-166.
802

A. New gene**B. Merged genes****C. Split****D. Changed orientation****E. Adjacent genes with ambiguous UTR coordinates****F. Introns with high read coverage and alternative splicing****G. Genes with overlapping coding sequences**

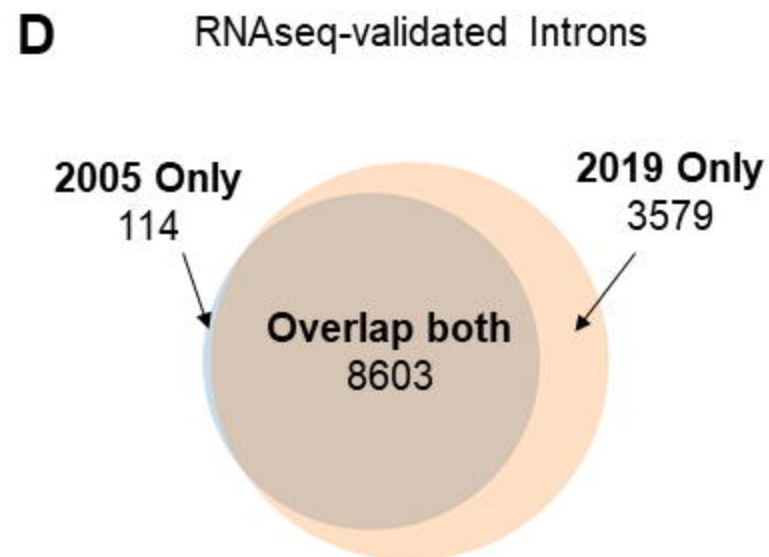
A

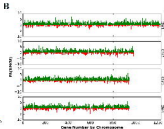
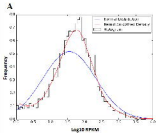
| Species | Proteome Size (#) | % ≥ 1 PFAM hit | Year Published |
|--------------------------------------|-------------------|---------------------|----------------|
| <i>Theileria parva</i> Muguga (2019) | 4055 | 61% | 2019 |
| <i>Theileria parva</i> Muguga (2005) | 4079 | 58% | 2005 |
| <i>Theileria annulata</i> Ankara | 3792 | 60% | 2005 |
| <i>Theileria equi</i> WA | 5332 | 53% | 2012 |
| <i>Theileria orientalis</i> Shintoku | 4002 | 59% | 2012 |
| <i>Babesia bovis</i> T2Bo | 3703 | 63% | 2007 |
| <i>Plasmodium vivax</i> Sall | 5393 | 62% | 2008 |
| <i>Plasmodium falciparum</i> 3D7 | 4555 | 61% | 2002 |

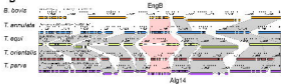


C

| 2005 | | 2019 | |
|--------------|---------------|--------------|---------------|
| Splice Sites | Frequency (#) | Splice Sites | Frequency (#) |
| GT...AG | 10,406 | GT...AG | 12,445 |
| GC...AG | 5 | GC...AG | 45 |
| TG...TA | 4 | | |
| TG...AA | 2 | | |
| AT...TG | 1 | | |
| AA...TA | 1 | | |
| GT...AA | 1 | | |
| AT...TA | 1 | | |
| AC...TA | 1 | | |
| Total | 10,422 | Total | 12,490 |





A*T. parva* Alg14**B****C**