

1 **The transcriptome of the avian malaria parasite *Plasmodium***  
2 ***ashfordi* displays host-specific gene expression**

3  
4  
5

6 **Authors**

7 Elin Videvall<sup>1</sup>, Charlie K. Cornwallis<sup>1</sup>, Dag Ahrén<sup>1,2</sup>, Vaidas Palinauskas<sup>3</sup>, Gediminas Valkiūnas<sup>3</sup>,  
8 Olof Hellgren<sup>1</sup>

9

10 **Affiliation**

11 <sup>1</sup>Department of Biology, Lund University, Sölvegatan 37, SE-22362 Lund, Sweden

12 <sup>2</sup>National Bioinformatics Infrastructure Sweden (NBIS), Lund University, Sölvegatan 37, SE-22362  
13 Lund, Sweden

14 <sup>3</sup>Institute of Ecology, Nature Research Centre, Akademijos 2, LT-08412 Vilnius, Lithuania

15

16 **Corresponding authors**

17 Elin Videvall ([elin.videvall@biol.lu.se](mailto:elin.videvall@biol.lu.se))

18 Olof Hellgren ([olof.hellgren@biol.lu.se](mailto:olof.hellgren@biol.lu.se))

19

20 **Running title**

21 The transcriptome of *Plasmodium ashfordi*

22

23 **Keywords**

24 *Plasmodium*, transcriptome, RNA-seq, malaria, gene expression, host-parasite interaction

25

26

## 27 **Abstract**

28

29 Malaria parasites (*Plasmodium* spp.) include some of the world's most widespread and virulent  
30 pathogens. Our knowledge of the molecular mechanisms these parasites use to invade and exploit  
31 hosts other than mice and primates is, however, extremely limited. It is therefore imperative to  
32 characterize transcriptome-wide gene expression from non-model malaria parasites and how this  
33 varies across host individuals. Here, we used high-throughput Illumina RNA-sequencing on blood  
34 from wild-caught Eurasian siskins experimentally infected with a clonal strain of the avian malaria  
35 parasite *Plasmodium ashfordi* (lineage GRW2). By using a multi-step approach to filter out host  
36 transcripts, we successfully assembled the blood-stage transcriptome of *P. ashfordi*. A total of 11 954  
37 expressed transcripts were identified, and 7 860 were annotated with protein information. We  
38 quantified gene expression levels of all parasite transcripts across three hosts during two infection  
39 stages – peak and decreasing parasitemia. Interestingly, parasites from the same host displayed  
40 remarkably similar expression profiles during different infection stages, but showed large differences  
41 across hosts, indicating that *P. ashfordi* may adjust its gene expression to specific host individuals.  
42 We further show that the majority of transcripts are most similar to the human parasite *Plasmodium*  
43 *falciparum*, and a large number of red blood cell invasion genes were discovered, suggesting  
44 evolutionary conserved invasion strategies between mammalian and avian *Plasmodium*. The  
45 transcriptome of *P. ashfordi* and its host-specific gene expression advances our understanding of  
46 *Plasmodium* plasticity and is a valuable resource as it allows for further studies analysing gene  
47 evolution and comparisons of parasite gene expression.

48

49

## 50 Introduction

51

52 The apicomplexan parasites of the genus *Plasmodium* (malaria parasites) encompass a worldwide  
53 distribution and infect a multitude of vertebrate hosts, including reptiles, birds and mammals  
54 (Garnham 1966). Their virulence can be highly variable between different strains and species. Some  
55 induce mild pathogenic effects on hosts and some cause severe disease, leading to high mortality rates  
56 (Palinauskas *et al.* 2008). Host individuals and host species also differ in their resistance and tolerance  
57 to malaria, and this interaction between host and parasite ultimately determines disease severity.  
58 Furthermore, the molecular response of hosts changes during the course of infection and creates a  
59 dynamic environment in which the parasites need to accommodate. Nevertheless, our understanding  
60 of how malaria parasites respond molecularly to different host individuals and to changes in the host  
61 immune defence over time is very limited.

62

63 Parasites of two clades in *Plasmodium* have been extensively studied from a molecular perspective,  
64 murine and primate parasites. We have learned a great deal about how malaria parasites of humans  
65 evolved and function by studying transcriptomes of their rodent-infecting relatives (Hall *et al.* 2005;  
66 Spence *et al.* 2013; Otto *et al.* 2014a). The majority of studies investigating gene expression in human  
67 malaria parasites have been conducted using cell lines (*in vitro*) or tissue cultures (*ex vivo*), which has  
68 provided tremendous insight into the biology of *Plasmodium* life stages (see e.g. Bozdech *et al.* 2003;  
69 Otto *et al.* 2010; Siegel *et al.* 2014). However, several discrepancies in parasite expression between  
70 cultures and live animals (*in vivo*) have been documented (Lapp *et al.* 2015) and a wide range of host  
71 environmental factors are absent in the *in vitro* systems. For example, temperature fluctuations,  
72 inflammatory and immune effector molecules, hormones, metabolites, microenvironments, and  
73 varying levels of oxygen, pH, and glucose are difficult to simulate in *in vitro* settings (LeRoux *et al.*  
74 2009). Parasites cultured outside hosts reflect this with different expression patterns, and markedly  
75 downregulate important vaccine candidate genes such as cell surface antigens (Daily *et al.* 2005; Siau  
76 *et al.* 2008). The natural host environment, which includes genotypic and immunological cues, may

77 therefore strongly affect the transcriptional responses of malaria parasites. To obtain representative  
78 transcriptional information from *Plasmodium* parasites, it is therefore valuable to study natural host  
79 systems.

80

81 The molecular mechanisms that enable successful invasion and establishment of malaria parasites in  
82 hosts other than mice and primates are unfortunately poorly known. Which genes are conserved  
83 across *Plasmodium* and how do virulence, immune evasion, and host-specificity vary in species  
84 infecting non-mammalian animals? To investigate these questions, it will be necessary to assemble  
85 and characterize genome-wide expression information from malaria parasites across their host  
86 phylogenetic range. With recent developments of high-throughput sequencing techniques, it has now  
87 become possible to generate genomic sequences from non-model parasites (Martinsen & Perkins  
88 2013). Dual RNA-sequencing of hosts and their parasites opens up possibilities of simultaneously  
89 studying host-parasite interactions and describing transcriptome expression in both actors.

90 Assembling novel parasite sequences *de novo* is, however, a very difficult task. Without a reference  
91 genome, transcripts from the host and/or other sources of contamination may remain after annotation  
92 and can influence downstream analyses. Meticulous filtering of assemblies using bioinformatics is  
93 therefore crucial to avoid erroneous conclusions (see e.g. Koutsovoulos *et al.* 2016). Nevertheless,  
94 successfully constructing parasite transcriptome data of high-quality from different hosts will provide  
95 us with valuable insights into the hidden biology of *Plasmodium*.

96

97 Malaria parasites that infect birds provide an excellent opportunity for studying transcriptional  
98 parasite responses due to their enormous diversity and large variation in host specificity and virulence  
99 (Bensch *et al.* 2004; Križanauskienė *et al.* 2006; Lachish *et al.* 2011). They are closely related to  
100 mammalian *Plasmodium*, but it was only recently established that rodent and primate malaria  
101 parasites are indeed monophyletic in the *Plasmodium* phylogeny (Bensch *et al.* 2016). Some avian  
102 *Plasmodium* are extreme host generalists, successfully infecting birds over several orders, while other  
103 parasites are host specialists infecting a single species (Pérez-Tris *et al.* 2007; Drovetski *et al.* 2014).  
104 The avian malaria system allows for the possibility of capturing wild birds in a natural setting, and

105 evaluating their status of malaria infection. Additionally, recent infection experiments have illustrated  
106 the potential to study *Plasmodium* in passerines under controlled conditions in laboratories  
107 (Zehntindjiev *et al.* 2008; Palinauskas *et al.* 2008, 2011; Cornet *et al.* 2014; Dimitrov *et al.* 2015; Ellis  
108 *et al.* 2015).

109 In this study, we used a bioinformatic multi-step filtering approach to assemble the blood  
110 transcriptome of the avian malaria parasite *Plasmodium ashfordi*, mitochondrial lineage GRW2. This  
111 parasite was first associated with various warblers (families Acrocephalidae and Phylloscopidae)  
112 (Valkiūnas *et al.* 2007), but has been found since 2016 in 15 different host species of two avian orders  
113 (Bensch *et al.* 2009). We used our assembly of *P. ashfordi* to evaluate transcriptome characteristics,  
114 genome-wide sequence similarity to other apicomplexans, and searched for genes known to be  
115 involved in the *Plasmodium* red blood cell invasion process. We analysed expression levels of  
116 parasite genes in three experimentally infected birds during two infection stages, peak and decreasing  
117 parasitemia, where the hosts have previously been shown to exhibit different transcriptome responses  
118 (Videvall *et al.* 2015). This allowed us for the first time to follow and describe the transcriptome of an  
119 avian malaria parasite over time in individual hosts.

120

## 121 **Results**

122

### 123 **The *Plasmodium ashfordi* transcriptome assembly**

124 We sequenced blood collected from three experimentally infected Eurasian siskins (*Carduelis spinus*)  
125 during peak and decreasing parasitemia (day 21 and 31 postinfection) with Illumina dual RNA-  
126 sequencing (see Methods for details). The transcriptome of *P. ashfordi* was assembled into two  
127 versions in order to make it transparent and as useful as possible for other researchers. The first  
128 assembly version, which we refer to as the annotated assembly, contains the transcripts with  
129 annotation information from proteins of Apicomplexa (n = 7 860) (Figure 1; Table S1, Supporting  
130 information). The second version that we refer to as the total assembly, also includes the unannotated  
131 contigs (n = 4 094) that were strictly filtered to remove contigs from the host (Figure 1), resulting in a  
132 total of 11 954 representative transcripts (Table 1). The genomes of *Plasmodium* parasites generally  
133 contain around 5-6 000 protein-coding genes (Kersey *et al.* 2016), making it reasonable to assume  
134 similar gene numbers for *P. ashfordi*. Building eukaryotic transcriptomes *de novo*, however, naturally  
135 yields many more transcripts than there are genes in the genome (Grabherr *et al.* 2011), and so the  
136 larger number of transcripts in our assembly is a result of isoform varieties, fragmented contigs, and  
137 non-coding RNA.

138         The size of the total transcriptome assembly is 9.0 Mbp and the annotated version of the  
139 assembly is 7.3 Mbp (81.20%) (Table 1). We calculated assembly statistics using the transcriptome-  
140 specific measurement E90N50, which represents the contig N50 value based on the set of transcripts  
141 representing 90% of the expression data, and is preferable over the original N50 when evaluating  
142 transcriptome assemblies (Haas 2016). The assembly E90N50 is 1 988 bp and the mean transcript  
143 length of the annotated transcripts is 930.8 bp (Table 1; Figure S1, Supporting information). In  
144 comparison, the length of coding sequences in the genome of *Haemoproteus tartakovskyi* (bird  
145 parasite in a sister genus to *Plasmodium*) is mean = 1 206 and median = 1 809 bp (Bensch *et al.*  
146 2016). The human parasite *Plasmodium falciparum* has a transcriptome with transcripts of median  
147 length 1 320 and mean length 2 197 bp (Gardner *et al.* 2002). The longest contig in the *P. ashfordi*

148 assembly, consisting of 26 773 bp, is transcribed from the extremely long ubiquitin transferase gene  
149 (AK88\_05171), which has a similar transcript length of around 27 400 bp in other *Plasmodium*  
150 species. The annotated transcriptome has an exceptionally low mean GC content of 21.22% (Figure  
151 1B), which is even lower than the highly AT-biased transcriptome of *P. falciparum* (23.80%).

152

### 153 **Functional analysis suggests that many host-interaction genes have evolved beyond recognition**

154 To evaluate biological and molecular functions of all annotated transcripts in *P. ashfordi*, we first  
155 analysed their associated gene ontology. An analysis of the transcriptome of *P. falciparum* was  
156 performed simultaneously to get an appreciation of how the *P. ashfordi* assembly compares  
157 functionally to a closely related, well-studied species. Overall, the two transcriptomes displayed  
158 highly similar gene ontology patterns (Figures 2A-B; Tables S8-S14, Supporting information).  
159 Transcripts primarily belonged to the two major molecular functions: ‘binding’ and ‘catalytic  
160 activity’, as well as the biological groups: ‘metabolic process’ and ‘cellular process’. These broad  
161 functional classes were also the gene ontology terms where *P. ashfordi* displayed a greater number of  
162 transcripts relative *P. falciparum* (Figure 2D). They are likely to contain multiple transcripts per gene  
163 due to isoforms, fragmented contigs, or possible gene duplications. The categories ‘receptor activity’,  
164 ‘cell adhesion’, and ‘multi-organism process’ were in contrast almost exclusively occupied by  
165 transcripts from *P. falciparum*. Interestingly, these categories predominately relate to the interaction  
166 with host cells and host defences (Figure 2C), which are known to contain genes showing strong  
167 evidence for positive selection in *Plasmodium* (Jeffares *et al.* 2007). It is highly likely that many *P.*  
168 *ashfordi* transcripts that belong to these host interaction processes have diverged sufficiently from  
169 mammalian parasite genes to escape annotation, but are expressed and present in the unannotated  
170 portion of the transcriptome assembly.

171 We explored the metabolic processes of *P. ashfordi* by looking into the child terms of this  
172 gene ontology. All metabolic categories of *P. ashfordi* contained similar or slightly more transcripts  
173 compared to *P. falciparum*, except the ‘catabolic process’ where *P. ashfordi* had fewer transcripts  
174 (Table S12, Supporting information). An investigation of the broad gene ontology category ‘kinase  
175 activity’ resulted in a total of 235 transcripts (Table S1, Supporting information). The kinase

176 multigene family FIKK has drastically expanded in the genomes of *P. falciparum* and *P. reichenowi*  
177 (Otto *et al.* 2014b), but only exists as one copy in the rodent malaria parasites. We used the proteins in  
178 the FIKK family from the genomes of *P. falciparum* (n = 21) and *P. yoelii* (n = 1) to search for  
179 matches in our assembly. All 22 protein sequences produced significant blast matches (e-value < 1e-  
180 5) against a single *P. ashfordi* transcript (TR71526|c0\_g2\_i1), further indicating that the FIKK gene  
181 expansion in primate malaria parasites likely happened after the split from both avian and rodent  
182 *Plasmodium*. Complete results of all gene ontology analyses of *P. ashfordi* and *P. falciparum* can be  
183 found in Tables S8-S14 (Supporting information).

184

### 185 **Gene expression is similar across different stages of infection**

186 Next, we analysed expression levels of the *P. ashfordi* transcripts within individual hosts across the  
187 two parasitemia stages. We accounted for differences in parasitemia levels between hosts and time  
188 points, and any variation in sequencing depth between samples, by normalizing overall expression  
189 values according to the DESeq method (Anders & Huber 2010). We found that the parasites displayed  
190 very similar gene expression patterns during peak and decreasing parasitemia stages (Figure 3A-E).  
191 No genes were significantly differentially expressed between the time points (q-value > 0.99), and the  
192 correlation in gene expression was extremely high (Pearson's product-moment correlation = 0.9983, t  
193 = 1 905.2, df = 11 952, p-value < 2.2e-16) (Figure 3A; Table S2, Supporting information). Annotated  
194 transcripts showing the highest expression fold change (non-significantly) between the two  
195 parasitemia stages were derived from the following genes (in order of most observed change): rho-  
196 GTPase-activating protein 1, 40S ribosomal protein S3a, two uncharacterized proteins, TATA-box-  
197 binding protein, heat shock protein 90, and C50 peptidase (Figure 3D; Table S2, Supporting  
198 information).

199

### 200 **Gene expression is host-specific**

201 In contrast to the similarities in gene expression between parasitemia stages, the parasite  
202 transcriptomes showed much larger differences in expression levels between the different host  
203 individuals. A principal component analysis of expression variance clustered parasite samples



204 together within their respective hosts, which showed major similarities in expression profiles (Figure  
205 3B). Samples derived from the same host individual did not separate until the third (15% variance)  
206 and fourth (13% variance) principal component dimensions (Figure 3C). The parasite transcriptome  
207 from host 4 during decreasing parasitemia showed the largest variation in parasite gene expression  
208 among all samples, yet it was still most similar to the transcriptome from the same host during peak  
209 parasitemia (Figure 3B; Figure 3E). In fact, all parasite transcriptomes during the decreasing  
210 parasitemia stage demonstrated closest distance to the transcriptome sample within the same host ten  
211 days earlier (Figure 3E). We further evaluated if specific transcripts contributed to the differences in  
212 parasite gene expression levels between individual hosts by performing a likelihood ratio test over all  
213 host individuals while controlling for parasitemia stage (time point). This resulted in 28 significant *P.*  
214 *ashfordi* transcripts (q-value < 0.1) displaying very high expression variation between hosts (Figure 4;  
215 Table S3, Supporting information). The most significant transcripts were derived from the genes  
216 cytochrome c oxidase subunit 1, 70 kd heat shock-like protein, M1 family aminopeptidase, and  
217 metabolite/drug transporter (Table S3, Supporting information).

218

### 219 ***P. ashfordi* is genetically identical in all hosts**

220 In an effort to investigate potential mechanisms behind the host-specific gene expression, we first  
221 evaluated if different haplotypes (multiclinality) of the parasite were present and had established  
222 differentially in the host individuals. As described in the Methods, all hosts were infected with the  
223 same clonal parasite isolate determined by the cytochrome b locus. This was also independently  
224 verified in all samples by examining read mapping of mitochondrial transcripts (Figure S5,  
225 Supporting information). Because sexual recombination of *Plasmodium* takes place in the mosquito,  
226 multiple alleles of nuclear genes could have been present in the parasite strain injected into the birds.  
227 A sensitive single nucleotide polymorphism (SNP) analysis over the entire *P. ashfordi* transcriptome  
228 found, however, extremely little variation in the parasite. During peak infection, we recovered a total  
229 of 10 (unfiltered) SNPs in the 11 954 *P. ashfordi* transcripts from the parasites in host 2, 32 SNPs in  
230 host 3, and 46 SNPs in host 4. The variation in number of SNPs called between parasites in different  
231 host individuals was due to read coverage, which is directly dependent on parasitemia levels, where

232 host 2 had the lowest, host 3 intermediate, and host 4 highest parasitemia (see Methods). After  
233 filtering on read depth, a total of 19 SNPs were identified and used in analyses (Table S4, Supporting  
234 information). We discovered that all SNPs (100%) were present in the parasites of all three host  
235 individuals, which carried the exact same allele polymorphism (e.g. C/T). To determine if the two  
236 alleles were expressed differentially in the hosts, allele frequency was calculated for each SNP in the  
237 parasite transcriptome from each host during peak infection. An analysis of the SNPs found no  
238 differences in allele frequencies between host individuals (Friedman rank sum test, Friedman chi-  
239 squared = 1.68, df = 2, p-value = 0.432), indicating similar allele expression levels in all parasite  
240 samples (Figure 3F). The SNPs were also found to be present in the parasites during the decreasing  
241 parasitemia stage, but lower coverage during this time point prevents statistical testing.

242 The 19 SNPs were located inside nine out of 11 954 contigs in the *P. ashfordi* transcriptome,  
243 resulting in an allelic occurrence of 0.075%. For a visual example of a contig with three SNPs present,  
244 see Figure S6 (Supporting information). Four of the nine transcripts containing SNPs were  
245 unannotated, and the others were derived from the genes merozoite surface protein 9 (MSP9), surface  
246 protein P113, ubiquitin related protein, multidrug resistance protein, and a conserved *Plasmodium*  
247 protein with unknown function (Table S4, Supporting information). We were also able to classify half  
248 of the SNPs whether they had any direct effects on the amino acid produced, and found that six SNPs  
249 were synonymous and three SNPs were non-synonymous. The three non-synonymous SNPs were  
250 located in transcripts from merozoite surface protein 9 (MSP9), ubiquitin related protein, and  
251 multidrug resistance protein (Table S4, Supporting information).

252 The presence of extremely few sequence polymorphisms in the transcriptome of *P. ashfordi*  
253 demonstrates that the parasite strain used was not only clonal with respect to the mitochondrial  
254 lineage, but homogeneous over the vast majority of the transcriptome. Furthermore, the detection of  
255 19 SNPs that were identical in the parasites from all host individuals verifies that the parasite is  
256 genetically identical in all hosts. Intriguingly, it also means that the *P. ashfordi* isolate originally  
257 consisted of a minimum of two different haplotypes which differed at nine genes before the  
258 experiment, and that this genetic variation managed to survive in the parasite population during the  
259 three week long infection, remaining in similar frequency in all hosts.

260

### 261 **Parasite sexual development does not differ between hosts**

262 We further evaluated the possibility that transcriptome host clustering were due to consistent  
263 differences relating to the sexual development of *P. ashfordi* during both time points in the host  
264 individuals. Lemieux *et al.* (2009) found that differences in *P. falciparum* gene expression between  
265 blood samples from children could be partially explained because the parasite expressed different  
266 genes during its sexual development (gametocyte) stage. We counted gametocytes in blood slides  
267 from the different samples using microscopy (Table S5, Supporting information) and found no  
268 differences in gametocyte proportions between host individuals and time points (Pearson's chi-  
269 squared = 1.59, df = 2, p-value = 0.452, mean peak infection = 6.74%, mean decreasing infection =  
270 10.24%). Furthermore, we directly compared the *P. ashfordi* transcripts showing significant  
271 expression differences between hosts (n = 28) to a list of sexual development genes exhibiting  
272 gametocyte-specific expression patterns in *P. falciparum* (n = 246) (Young *et al.* 2005). Only one of  
273 the 28 genes had a match in the sexual development gene set, namely the metabolite/drug transporter  
274 MFS1 (Table S3, Supporting information). However, this single match is not greater than expected by  
275 chance (hypergeometric test, p-value = 0.372), given probability of a match across all protein-coding  
276 genes in the *P. falciparum* genome (n = 5 344). Together, these results show that the host-specific  
277 expression pattern in *P. ashfordi* cannot be explained due to differences in sexual development.

278

### 279 ***P. ashfordi* shows sequence similarities to human malaria parasites**

280 Almost all annotated contigs (99.59%; n = 7 828) resulted in a best blast hit against a species within  
281 the genus *Plasmodium* (Figure 5A). The remaining contigs had matches against species within the  
282 genera of *Eimeria* (n = 12), *Cryptosporidium* (n = 6), *Neospora* (n = 5), *Babesia* (n = 4), *Hammondia*  
283 (n = 2), *Ascogregarina* (n = 1), *Theileria* (n = 1), and *Toxoplasma* (n = 1) (Table S6, Supporting  
284 information). The great majority (73.59%) of the contig blast matches were proteins originating from  
285 primate parasites, while 25.34% matched rodent parasites, and only 0.92% parasites of birds (Figure  
286 5B).

287 At the species level, most contigs (29.91%) resulted in best blast hit against *P. falciparum*,  
288 followed by *P. reichenowi* (16.88%) and *P. yoelii* (8.59%) (Figure 5C). The significant blast matches  
289 to bird parasites consisted of the species *Plasmodium gallinaceum* (n = 56), *Eimeria acervulina* (n =  
290 5), *Eimeria tenella* (n = 4), *Eimeria mitis* (n = 3), *Plasmodium relictum* (n = 3), and *Plasmodium lutzii*  
291 (n = 1). The contigs giving matches to avian *Plasmodium* were primarily derived from commonly  
292 sequenced apicomplexan genes and therefore available in public databases, for example cytochrome c  
293 oxidase subunit 1 (COX1; *P. lutzii*), merozoite surface protein 1 (MSP1; *P. relictum*), thrombospondin  
294 related anonymous protein (TRAP; *P. relictum*), and cytochrome b (CYTB; *P. gallinaceum*) (Table  
295 S1, Supporting information).

296 The five contigs with highest GC content in the *P. ashfordi* transcriptome (47.7% – 56.4%)  
297 all had matches against the avian parasites *Eimeria*, despite them only comprising 0.15% (n = 12) of  
298 the total annotation. *Eimeria* species have a very high transcriptome GC content (*E. acervulina*:  
299 55.98%; *E. mitis*: 57.30%), and the *P. ashfordi* transcripts matching this genus consist mostly of  
300 ribosomal and transporter genes (Table S1, Supporting information). The *P. ashfordi* contigs with  
301 highest expression levels were primarily annotated by uncharacterized protein matches to the rodent  
302 parasite *P. yoelii* (Table S7, Supporting information). In fact, the six most highly expressed transcripts  
303 that were annotated, all gave significant blast matches to *P. yoelii*. Further investigation revealed that  
304 these transcripts are most likely derived from ribosomal RNA.

305

### 306 **Identification of conserved *Plasmodium* invasion genes**

307 Finally, to assess molecularly conserved strategies of *P. ashfordi* compared to mammalian malaria  
308 parasites, we searched for annotated genes known to be involved in the red blood cell invasion by  
309 *Plasmodium*. (Bozdech *et al.* 2003; Beeson *et al.* 2016). We discovered successfully assembled *P.*  
310 *ashfordi* transcripts from a whole suite of host cell invasion genes (Table 2). This includes for  
311 example the genes merozoite surface protein 1 (MSP1), apical membrane antigen 1 (AMA1),  
312 merozoite adhesive erythrocytic binding protein (MAEBL), GPI-anchored micronemal antigen  
313 (GAMA), and the rhoptry neck binding proteins 2, 4, and 5 (RON2, RON4, and RON5). Interestingly,  
314 the *P. ashfordi* RON genes in particular seemed to slightly decrease expression levels in all hosts over

315 the two time points (Figure 6). In general, however, the invasion genes showed a range of expression  
316 patterns over time, going in various directions (Figure 6).

317 All genes known to be involved in the *Plasmodium* motor complex (Opitz & Soldati 2002;  
318 Baum *et al.* 2006), driving parasite gliding motion and enabling host cell invasion, were discovered in  
319 *P. ashfordi*. These include: actin (ACT1), actin-like protein (ALP1), aldolase (FBPA), myosin A  
320 (MyoA), myosin A tail interacting protein (MTIP), glideosome-associated protein 45 and 50 (GAP45  
321 and GAP50), and thrombospondin related anonymous protein (TRAP) (Table 2). We also found the  
322 bromodomain protein 1 (BDP1), which has been directly linked to erythrocyte invasion by binding to  
323 chromatin at transcriptional start sites of invasion-related genes and controlling their expression  
324 (Josling *et al.* 2015).

325 We found two transcripts matching the low molecular weight rhoptry-associated proteins 1  
326 and 3 (RAP1 and RAP3) that are secreted from the rhoptry organelles during cell invasion. The  
327 genomes of human malaria parasites contain a paralog gene called RAP2 as well, whereas rodent  
328 malaria parasites contain a single gene copy that is a chimera of RAP2 and RAP3 (RAP2/3)  
329 (Counihan *et al.* 2013). The *P. ashfordi* transcript in question (TR13305|c0\_g1\_i1) matches *P.*  
330 *falciparum* RAP3 better than the rodent parasite version of RAP2/3. The three high molecular weight  
331 rhoptry proteins (RhopH1, RhopH2, RhopH3) which bind to the erythrocyte plasma membrane and  
332 transfer to the parasitophorous vacuole membrane upon invasion (Vincensini *et al.* 2008; Counihan *et*  
333 *al.* 2013) were all identified in *P. ashfordi*. RhopH1 encompasses the multigene family of  
334 cytoadherence linked asexual proteins (CLAGs), present in varying copy number across *Plasmodium*.

335 Other assembled *P. ashfordi* orthologs of genes involved in host cell invasion were the  
336 rhoptry-associated leucine zipper-like protein 1 (RALP1), rhoptry-associated membrane antigen  
337 (RAMA), armadillo-domain containing rhoptry protein (ARO), RH5 interacting protein (RIPR),  
338 TRAP-like protein (TLP), merozoite TRAP-like protein (MTRAP), thrombospondin related apical  
339 membrane protein (TRAMP), subtilisin proteases 1 and 2, (SUB1 and SUB2), and merozoite surface  
340 proteins 8 and 9 (MSP8 and MSP9). MTRAP and TRAMP are proteins that belong to the TRAP-  
341 family and are released from the microneme organelles during invasion (Green *et al.* 2006; Cowman  
342 *et al.* 2012), and the subtilisin proteases SUB1 and SUB2 are heavily involved in the processing and

343 cleavage of immature merozoite antigens, for example MSP1 and AMA1 (Beeson *et al.* 2016). ARO  
344 plays a crucial role in positioning the rhoptry organelles within the apical end of the parasite to enable  
345 the release of rhoptry-associated molecules (Mueller *et al.* 2013) such as RAMA and RALP1, which  
346 then bind to the erythrocyte surface.

347         We furthermore discovered transcripts of several reticulocyte binding proteins (RBP/RH)  
348 thought to be absent in the genomes of avian malaria parasites (Lauron *et al.* 2015). These particular  
349 transcripts, together with RAMA, showed much higher e-values than other invasion genes (Table 2),  
350 indicating high differentiation between avian and mammalian *Plasmodium* RH genes. Finally, two  
351 rhomboid proteases (ROM1 and ROM4) have been linked to host cell invasion in *P. falciparum* via  
352 cleavage of transmembrane adhesins (Baker *et al.* 2006; Santos *et al.* 2012). We found both of these  
353 genes, together with other rhomboid proteases (ROM2, ROM3, ROM6, ROM8, and ROM10)  
354 expressed in *P. ashfordi*. More information about the assembled genes can be found in Table S1  
355 (Supporting information).

356

357

## 358 **Discussion**

359

360 In this study, we assembled and characterized a blood-stage transcriptome with quantified gene  
361 expression of an avian malaria parasite, *P. ashfordi*. By developing a bioinformatic filtering method  
362 capable of dealing with dual RNA-seq data, we effectively removed contigs originating from the host  
363 and other sources of contamination in a multistep approach (Figure 1). This resulted in a  
364 transcriptome with 7 860 annotated transcripts and an additional 4 094 unannotated transcripts. We  
365 discovered that 19 SNPs were not only present, but identical, in the transcriptomes of all six parasite  
366 samples from the three hosts, which corroborates that these are indeed true sequence polymorphisms  
367 and is an excellent independent verification that all host transcripts have been successfully removed  
368 from the assembly (Table S4, Supporting information). The gene expression of *P. ashfordi* displayed  
369 strikingly similar patterns during peak and decreasing infection stages and within individual hosts  
370 (Figure 3). Furthermore, *P. ashfordi* shows most sequence similarities to the human malaria parasite  
371 *P. falciparum* (Figure 5C), but specific genes involved in host interaction and defence seem to be  
372 highly differentiated between the two parasites (Figure 2C). Nonetheless, the assembly supports  
373 several important erythrocyte invasion genes (Table 2), indicating evolutionary conserved cell  
374 invasion strategies across the phylogenetic host range of *Plasmodium* parasites.

375

### 376 ***P. ashfordi* displays host-specific gene expression**

377 Interestingly, and contrary to our expectations, *P. ashfordi* showed highly similar expression profiles  
378 inside the same host, despite being sampled ten days apart during two different disease stages. All  
379 birds were inoculated with the same malaria strain derived from a single donor bird, and our SNP  
380 analysis verified that the genetic composition of *P. ashfordi* is identical in all hosts. The mechanism  
381 behind this host-specific expression pattern is unknown, but there is no reason why six independently  
382 sampled parasite populations should cluster transcriptionally based on host individual unless  
383 expression levels are somehow regulated in response to hosts. The expression pattern can potentially  
384 be caused by genotype by genotype interactions between the host and the parasite, modulation of

385 parasite expression by the host, or plasticity of the parasite to different host environments. This result  
386 has potentially important implications for our understanding of the evolution of host-parasite  
387 interactions, and as a result warrants further research extending the limited sample size to more hosts  
388 and more timepoints throughout the infection.

389 Host genotype by parasite genotype interactions are complicated and not that well  
390 documented in malaria parasite systems. Studies with different genotypes of both host and parasite  
391 have found effects of host genotype, but not parasite genotype, on factors such as host resistance and  
392 parasite virulence (Mackinnon *et al.* 2002; de Roode *et al.* 2004; Grech *et al.* 2006; see also  
393 Idaghdour *et al.* 2012). Less is known about the transcriptome responses of malaria parasites to  
394 different host individuals. Some studies have found differential gene expression responses of  
395 *Plasmodium* to resistant versus susceptible mice strains (see e.g. Lovegrove *et al.* 2006), and Daily *et*  
396 *al.* (2007) discovered host-specific distinct transcriptional states of *P. falciparum* in the blood of  
397 Senegalese children. However, a reanalysis of the data by Daily *et al.* (2007) suggested that  
398 differences in sexual development of the parasite may have contributed to the transcriptional states  
399 (Lemieux *et al.* 2009). With respect to our results, we found no evidence of differences in sexual  
400 development of the parasite and our SNP analysis showed that parasite haplotypes did not establish  
401 differentially in the host individuals. This suggests that sexual development and haplotype differences  
402 had little influence on our results and that the most likely explanation for the host-specific  
403 transcriptome patterns is likely to be plasticity in gene expression by *P. ashfordi*.

404 The expression profiles of *P. ashfordi* did not exhibit any significant differences between  
405 peak and decreasing parasitemia stages (day 21 and 31 post-infection). The hosts in our study  
406 experienced relatively high parasitemia levels during the decreasing parasitemia stage as well (see  
407 Methods), so it is possible that these specific time points do not provide very different environmental  
408 host cues to the parasites. However, the concurrent transcriptomes of the avian hosts (analysed in  
409 Videvall *et al.* 2015) displayed large differences in gene expression between these two parasitemia  
410 stages, notably with a reduced immune response during decreasing parasitemia. This is important  
411 because it appears that *P. ashfordi* does not adjust its gene expression in response to the decreasing  
412 immune levels of the hosts, but instead conforms to the specific environment of individual hosts.



413 *P. falciparum* evades the human immune defence via intracellularity, clonal antigenic  
414 variation (Guizetti & Scherf 2013), transcriptional antigenic switches (Recker *et al.* 2011), splenic  
415 avoidance by erythrocytic adherence to the endothelium (Craig & Scherf 2001) and sequestration in  
416 organ microvasculature (Silamut *et al.* 1999), erythrocytic rosetting (Niang *et al.* 2014), host  
417 immunosuppression (Hisaeda *et al.* 2004), and manipulation of host gene expression. It is possible  
418 that *P. ashfordi* shares several of these evasion strategies with *P. falciparum*, although this remains  
419 unknown. One example of immune evasion by manipulation of host gene expression is the parasite  
420 gene macrophage migration inhibitory factor homologue (MIF), which contributes to *Plasmodium*  
421 parasites' ability to modulate the host immune response by molecular mimicry (Cordery *et al.* 2007).  
422 This gene was discovered transcribed in *P. ashfordi* as well (TR2046|c0\_g1\_i1), suggesting that a  
423 similar immune manipulation strategy is possible (Table S1, Supporting information).

424

#### 425 **Similarities to *P. falciparum* and other malaria parasites**

426 The majority of all annotated contigs (73.59%) resulted in a best blast hit against primate parasites  
427 (Figure 5B). Curiously, the human malaria parasite *P. falciparum* comprised the majority of all  
428 matches with almost a third of the transcripts (29.91%) (Figure 5C). This is likely because *P.*  
429 *falciparum* currently constitutes the organism with most sequence similarities to *P. ashfordi* based on  
430 publically available sequences (Bensch *et al.* 2016). The chimpanzee parasite *P. reichenowi* had the  
431 second most blast matches to *P. ashfordi* (Figure 5C), and it is the closest living relative to *P.*  
432 *falciparum* based on current genomic data (Otto *et al.* 2014b). Furthermore, both *P. falciparum* and *P.*  
433 *ashfordi* share the genome characteristics of being extremely AT-biased, with *P. ashfordi* reaching a  
434 remarkably low transcriptomic GC content of 21.22% (Table 1) compared to the already AT-rich *P.*  
435 *falciparum*, which has a transcriptomic GC content of 23.80%. Lastly, because of its role in human  
436 disease, *P. falciparum* is the most sequenced *Plasmodium* species (172 official strains in the NCBI  
437 taxonomy database as of May 2016) (Gardner *et al.* 2002), resulting in the greatest opportunity for  
438 transcript sequences to find significant blast matches.

439 Less than one percent of all contigs resulted in a blast hit against avian parasites (0.92%).

440 This is due to the fact that almost no genomic resources are available for avian *Plasmodium*. Despite

441 their enormous diversity, world-wide distribution, and harmful effects on susceptible bird populations,  
442 genomic studies of avian malaria parasites have been largely non-existent until now. The genome of  
443 *P. gallinaceum*, the malaria parasite of chickens, has been sequenced but not published and we were  
444 therefore not able to use it in our analyses. A transcriptome assembly of *P. gallinaceum* is available  
445 for download (Lauron *et al.* 2014, 2015), although still contains a large proportion of contigs  
446 matching birds, making comparisons with *P. ashfordi* difficult (see Figure S4, Supporting  
447 information). Dual RNA-sequencing of a more distantly related apicomplexan parasite,  
448 *Leucocytozoon buteonis*, and its buzzard host was recently described (Pauli *et al.* 2015), though no  
449 publically available transcriptome exists. Finally, both 454 RNA-sequencing and Illumina genome  
450 sequencing of the generalist avian malaria parasite *P. relictum* (lineage SGS1) have been performed  
451 (Hellgren *et al.* 2013; Bensch *et al.* 2014; Lutz *et al.* 2016), but the extremely low sequence coverage  
452 in both cases does unfortunately not allow for assembly nor any genomic analyses. We hope that  
453 future sequencing of these avian parasites will enable genome-wide comparisons.

454         The lack of genome-wide sequence data from malaria parasites of hosts other than mice and  
455 primates means that little is known about which genes across *Plasmodium* are conserved and which  
456 that are unique. Our gene ontology results confirm that the transcriptomes of *P. ashfordi* and *P.*  
457 *falciparum* overall are functionally similar (Figure 2A-B), though specific genes involved in host  
458 interaction and receptor binding could not be directly located in the *P. ashfordi* assembly (Figure 2C).  
459 The transcriptome of *P. falciparum* is certainly more complete because it is based on the genome  
460 sequence (Gardner *et al.* 2002) and therefore includes genes expressed from the entire life cycle. As a  
461 result, any *P. falciparum* genes not found in *P. ashfordi*, are either specifically transcribed during  
462 certain life stages (e.g. in the mosquito), not present in the genome, or too diverged to be detected  
463 with sequence similarity searches. Seeing how this particular group of genes involved in host  
464 interaction are under strong positive and diversifying selection pressure in *Plasmodium* (Hall *et al.*  
465 2005; Jeffares *et al.* 2007; Otto *et al.* 2014b), it is likely that many are indeed present in the *P.*  
466 *ashfordi* assembly, but unannotated due to evolutionary divergence.

467         As a step to investigate genes involved in host interaction in more detail, we searched in the  
468 *P. ashfordi* assembly for genes specifically known to be involved in the merozoite invasion of host

469 red blood cells. Previously, only a handful of studies have sequenced candidate invasion genes in  
470 avian malaria parasites; these include MAEBL (Martinez *et al.* 2013), AMA1 and RON2 (Lauron *et*  
471 *al.* 2014), MSP1 (Hellgren *et al.* 2013, 2015), RIPR (Lauron *et al.* 2015), and TRAP (Templeton &  
472 Kaslow 1997; Farias *et al.* 2012). Due to the evolutionary distance between mammals and saurians,  
473 and their inherent blood cell differences (birds/reptiles have erythrocytes with a nucleus and  
474 mitochondria while mammalian cells are anucleated), we might expect to find few and highly  
475 differentiated gene orthologs. Instead, we discovered a large number of red blood cell invasion genes  
476 expressed in *P. ashfordi* (Table 2), indicating that most of these specific invasion genes are conserved  
477 across both mammalian and avian *Plasmodium*.

478         The invasion genes that were most differentiated between birds and mammals were the  
479 rhoptry-associated membrane antigen (RAMA) and the reticulocyte binding proteins (RBP/RH),  
480 which had diverged almost beyond recognition. These RH genes, together with other erythrocyte  
481 binding antigens (EBA), have been assumed to be absent in the genomes of avian malaria parasites  
482 (Lauron *et al.* 2015). However, our result suggests that several are not only present, but also  
483 transcribed, though with high sequence divergence. It is possible that additional erythrocyte binding  
484 proteins are present in the *P. ashfordi* assembly, though ortholog searches for these genes will become  
485 complicated if they have evolved under especially strong selection pressure in avian *Plasmodium*.

486

487

## 488 **Conclusion**

489 In this study we have *de novo* assembled, characterized, and evaluated the blood transcriptome of the  
490 avian parasite *Plasmodium ashfordi*. By developing a rigorous bioinformatic multistep approach, the  
491 assembly was successfully cleaned of host sequences and contains high numbers of important genes  
492 in e.g. red blood cell invasion. We have shown that *P. ashfordi* displays similar expression profiles  
493 within individual hosts during two different stages of the infection – but different expression patterns  
494 between individual hosts – indicating possible host specific parasite gene regulation. The expression  
495 information of all transcripts will assist researchers studying genes involved in e.g. immune evasion,

496 host-specificity, and parasite plasticity. In addition, our results show that the isolate of *P. ashfordi*  
497 used in this experiment originally contained two alleles at nine loci and has managed to maintain this  
498 low-level genomic heterozygosity throughout the infection in all host individuals. The results  
499 presented here and the associated assembly will help improve our understanding of host-parasite  
500 interactions, evolutionary conserved *Plasmodium* strategies, and the phylogenetic relationships  
501 between apicomplexans.  
502

## 503 **Methods**

504

### 505 **Experimental setup**

506 We used four wild-caught juvenile Eurasian siskins (*Carduelis spinus*) in an infection experiment.

507 The experimental procedure was carried out 2012 at the Biological Station of the Zoological Institute  
508 of the Russian Academy of Sciences on the Curonian Spit in the Baltic Sea (55° 05'N, 20° 44'E). All  
509 details regarding the setup have been outlined in Videvall *et al.* (2015). Three of the birds were  
510 inoculated with a single injection of blood containing the erythrocytic stages of *Plasmodium ashfordi*,  
511 subgenus *Novyella*, lineage GRW2. For a description of this parasite, see Valkiūnas *et al.* (2007). A  
512 control bird (bird 1) was simultaneously inoculated with uninfected blood for the evaluation of host  
513 transcription (see Videvall *et al.* 2015).

514 The parasite strain was originally collected in 2011 from a single common cuckoo (*Cuculus*  
515 *canorus*) that had acquired the infection naturally. It was thereafter multiplied in common crossbills  
516 (*Loxia curvirostra*) in the laboratory and deep frozen in liquid nitrogen for storage. One crossbill was  
517 subsequently infected with the parasite, and blood from this single bird was used to infect our  
518 experimental birds. Crossbills were used as intermediate hosts because they are both susceptible to  
519 this strain and have a large enough body size to provide the amount of blood needed for donations. A  
520 subinoculation of a freshly prepared mixture containing infected blood from the donor was made into  
521 the pectoral muscle of the recipient birds (details of procedure can be found in Palinauskas *et al.*  
522 2008). By using a single donor, we ensured that the same clonal parasite strain and parasite quantity  
523 was injected into recipient birds.

524 All birds were thoroughly screened with both microscopic (Palinauskas *et al.* 2008) and  
525 molecular (Hellgren *et al.* 2004) methods before the experiment to make sure they had no prior  
526 haemosporidian infection. Blood samples for RNA-sequencing were taken from birds before infection  
527 (day 0), during peak parasitemia (day 21 postinfection) and during the decreasing parasitemia stage  
528 (day 31 postinfection). The parasitemia intensity varied substantially in infected birds, with bird 2,  
529 bird 3, and bird 4 having 24.0%, 48.0%, and 71.3% of their red blood cells infected during peak

530 parasitemia, and later 8.2%, 21.8%, and 33.3%, respectively, during the decreasing parasitemia stage  
531 (Videvall *et al.* 2015). The two parasitemia stages are referred to as different ‘infection stages’ in the  
532 hosts, but we want to clarify that there is no evidence present suggesting that the parasites have  
533 entered a different stage of their life cycle (e.g. tissue merogony). Experimental procedures were  
534 approved by the International Research Co-operation Agreement between the Biological Station  
535 Rybachy of the Zoological Institute of the Russian Academy of Sciences and Institute of Ecology of  
536 Nature Research Centre (25 May 2010). Substantial efforts were made to minimize handling time and  
537 potential suffering of birds.

538

### 539 **RNA extraction and sequencing**

540 From six infected samples (three treatment birds at days 21 and 31) and six uninfected samples (three  
541 treatment birds at day 0 and one control bird at days 0, 21, and 31), total RNA was extracted from 20  
542  $\mu$ l whole blood. Detailed extraction procedures can be found in Videvall *et al.* (2015). Total extracted  
543 RNA was sent to Beijing Genomics Institute (BGI), China, for RNA quality control, DNase  
544 treatment, rRNA depletion, and amplification using the SMARTer Ultra Low kit (Clontech  
545 Laboratories, Inc.). BGI performed library preparation, cDNA synthesis, and paired-end RNA-  
546 sequencing using Illumina HiSeq 2000. The blood samples from bird 3 and bird 4 during peak  
547 parasitemia were sequenced by BGI an additional time in order to generate more reads from the  
548 parasite in preparation for this transcriptome assembly. These resequenced samples were regarded and  
549 handled as technical replicates. We quality-screened all the demultiplexed RNA-seq reads using  
550 FastQC (v. 0.10.1) (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>).

551

### 552 ***De novo* transcriptome assembly**

553 Quality-filtered RNA-seq reads from all six infected bird samples together with the two re-sequenced  
554 samples were used in a *de novo* assembly. This was performed using the transcriptome assembler  
555 Trinity (v. 2.0.6) (Grabherr *et al.* 2011) with 489 million reads. Mapping of reads to available  
556 genomes of human malaria parasites were employed simultaneously, but unfortunately these attempts

557 yielded few hits due to the evolutionary distance between avian *Plasmodium* and human *Plasmodium*,  
558 so we continued annotating the assembly using blast searches.

559         The assembled transcripts were blasted against the NCBI non-redundant protein database  
560 using the program DIAMOND BLASTX (v. 0.7.9) (Altschul *et al.* 1990; Buchfink *et al.* 2014) with  
561 sensitive alignment mode. A total of 47 823 contigs generated significant hits against avian species  
562 (Figure 1C). A large number of contigs (n = 260 162) did not produce any significant (e-value < 1e-5)  
563 blastx hits (Figure 1D). This is because 1) the host species is a non-model organism without a genome  
564 available, leading to a large number of host contigs without blast hits, 2) contigs might not necessarily  
565 contain coding sequences, but can be derived from noncoding RNAs, etc. and will therefore not match  
566 protein sequences, 3) short, fragmented contigs may not yield sufficient blast hit significance, and 4)  
567 an extreme underrepresentation of protein sequences from avian *Plasmodium* species in the NCBI nr  
568 database will not result in any significant blast hits to genes unique in avian malaria parasites.

569         We strictly filtered the initial assembly by only retaining a total of 9 015 transcripts  
570 (isoforms) that produced significant blast matches against proteins from species in the Apicomplexa  
571 phylum. A previous assembly using an earlier version of Trinity (v. r20140413p1) (Grabherr *et al.*  
572 2011) performed better when it came to assembling the longest contigs (> 6 kbp). Different versions  
573 of assembler software may construct de Bruijn graphs somewhat differently, which is why it can be a  
574 good idea to make several assemblies and later combine parts of them (Brian Haas, personal  
575 communication). The previous assembly had been blasted and screened for Apicomplexa in exactly  
576 the same way as described above. In order not to lose these important full-length transcripts, we  
577 therefore included the longest contigs from the previous assembly that had 1) not assembled correctly  
578 in the current assembly, and 2) significant blastx hits against Apicomplexa (n = 10), resulting in a  
579 total of 9 025 transcripts. The fact that these contigs contained similar sequences already present in  
580 the assembly was dealt with through downstream clustering of the sequences.

581

## 582 **Transcriptome cleaning and filtering**

583 Some contigs in the annotated assembly contained poly-A tails which complicated downstream  
584 analyses and resulted in biased mapping estimates. We therefore removed all poly-A/T tails using

585 PRINSEQ (v. 0.20.4) (Schmieder & Edwards 2011) with a minimum prerequisite of 12 continuous  
586 A's or T's in the 5' or 3' end of the contigs. A total of 106 202 bases (1.18%) was trimmed and the  
587 mean transcript length was reduced from 995.28 to 983.52 bases.

588         The unknown transcripts that failed to produce significant hits to any organism during the  
589 blastx run (n = 260 162) were subsequently cleaned using the following procedure. First, we trimmed  
590 them for poly-A tails, resulting in a total of 455 331 bases removed, and a slight decrease of the mean  
591 length of the unknown sequences from 555.27 nt before trimming to 553.52 nt after trimming. The  
592 majority of these unknown transcripts came from host mRNA, but their GC content displayed a clear  
593 bimodal distribution (Figure 1D), where the contigs with very low GC were strongly suspected to  
594 originate from the parasite. To avoid any host contigs, we strictly filtered the unknown transcripts to  
595 only include sequences with a mean GC content lower than 23% (n = 4 624). This threshold was  
596 based on the Apicomplexa-matching transcripts (mean GC = 21.22%), the contigs matching birds  
597 (class: Aves; n = 47 823; mean GC = 47.65%), and the bird contig with the absolute lowest GC  
598 content (GC = 23.48%) (Figure 1C).

599

#### 600 **Transcriptome clustering, further filtering, and validation**

601 To reduce redundancy of multiple isoforms and transcripts derived from the same gene, we first  
602 merged together the annotated and the unknown transcripts with a GC content < 23% (n = 13 649).  
603 We then clustered these sequences together in order to retain most transcripts but group the highly  
604 similar ones based on 97% sequence similarity and a k-mer size of 10, using CD-HIT-EST (v. 4.6) (Li  
605 & Godzik 2006). The most representative (longest) contig in every cluster was selected to represent  
606 those transcripts, resulting in 12 266 contigs/clusters. We further filtered all the short sequences (<  
607 200 bases), to obtain a set of 12 182 representative transcripts.

608         A second blast filtering step with the trimmed representative contigs against the Refseq  
609 genomic database was then employed using BLASTN+ (v. 2.2.29) (Altschul *et al.* 1990; Camacho *et*  
610 *al.* 2009) to identify some ambiguous contigs suspected to contain non-coding RNA bird sequences.  
611 We removed all contigs that gave significant matches (e-value < 1e-6) against all animals (kingdom:



612 Metazoa), so we could be confident that the assembly only consisted of true parasite transcripts. This  
613 last filtering step removed 228 contigs.

614 The unannotated transcripts (n = 4 094) of the final assembly were further validated to  
615 originate from the parasite by using reads from the six uninfected samples of the same hosts sampled  
616 before infection and the control bird. A total of 350 318 482 reads (65 bp and 90 bp) from all  
617 uninfected samples were mapped to the unannotated transcripts using Bowtie2 (v. 2.2.5) (Langmead  
618 & Salzberg 2012), resulting in the alignment of 90 read pairs (0.000051%). This extremely low  
619 mapping percentage from the uninfected samples greatly supported our conclusion that these  
620 transcripts had indeed been transcribed by *P. ashfordi*. These 4 094 representative transcripts with  
621 unknown function are referred to throughout the paper as unannotated transcripts. The resulting final  
622 transcriptome assembly consists of 11 954 representative annotated and unannotated transcripts.

623

#### 624 **Estimating expression levels**

625 Poly-A tails of 489 million RNA-seq reads from all samples were trimmed as well using PRINSEQ  
626 (v. 0.20.4) (Schmieder & Edwards 2011). A minimum prerequisite for trimming were 20 continuous  
627 A's or T's in the 5' or 3' end of each read. Only trimmed reads longer than 40 bp and still in a pair  
628 were retained (n = 451 684 626) in order to confidently map high quality reads with good minimum  
629 lengths. Bowtie2 (v. 2.2.5) (Langmead & Salzberg 2012) was used to map the trimmed RNA-seq  
630 reads of every sample (n = 8) (six biological and two technical replicates) back to the *P. ashfordi*  
631 transcriptome consisting of the 11 954 representative sequences. We calculated expression levels  
632 using RSEM (v. 1.2.21) (Li & Dewey 2011), which produces expected read counts for every contig.

633 The counts of the 11 954 transcripts were subsequently analysed inside the R statistical  
634 environment (v. 3.2.5) (R Core Team 2015). We tested for expression differences in the malaria  
635 parasites between the two time points, and between the hosts, using DESeq2 (v. 1.10.1) (Love *et al.*  
636 2014). The two resequenced samples (technical replicates) of bird 3 and bird 4 during peak  
637 parasitemia were handled exactly the same as all other samples, and their respective read count were  
638 added to their biological samples, according to the DESeq2 manual. Counts were normalized to  
639 account for potential variation in sequencing depth as well as the large differences in number of

640 parasites present in the blood (parasitemia levels). Regularized log transformation of counts was  
641 performed in order to represent the data without any prior knowledge of sampling design in the  
642 principal component analysis and sample distance calculations. This way of presenting counts without  
643 bias is preferred over variance stabilizing of counts when normalization size factors varies greatly  
644 between the samples (Love *et al.* 2014), as they naturally do in our data.

645

#### 646 **Variant calling and SNP analyses**

647 Sequence variation in the transcriptome assembly was performed according to the Genome Analysis  
648 Toolkit (GATK) Best Practices workflow for variant calling in RNA-seq data (v. 2015-12-07). The  
649 BAM files for each sample produced by RSEM were sorted, read groups added, and duplicates reads  
650 removed using Picard Tools (v. 1.76) (<https://broadinstitute.github.io/picard/>). With the  
651 SplitNCigarReads tool in GATK (v. 3.4-46) (McKenna *et al.* 2010), the mapped reads were further  
652 reassigned mapping qualities, cleaned of Ns, and hard-clipped of overhang regions. Variant calling  
653 was done in GATK with HaplotypeCaller using ploidy = 1 and optimized parameters recommended  
654 for RNA-seq data (see GATK Best Practices). Indels were excluded using the SelectVariants tool, and  
655 variants were filtered on quality using the VariantFiltration tool in GATK according to the filter  
656 recommendation for RNA-seq data. Next, we used SnpSift in SnpEff (v. 4\_3g) (Cingolani *et al.* 2012)  
657 to filter variants on depths of a minimum of 20 high-quality, non-duplicated reads. Finally, to  
658 calculate allele frequencies of all samples in variant positions called in only some hosts, we ran  
659 HaplotypeCaller again in the `-ERC BP_RESOLUTION` mode and extracted the read depths of the  
660 nucleotide positions called previously. Nucleotide positions were filtered if they were positioned in  
661 the very end of contigs or had a depth of less than 20 high-quality, non-duplicated reads according to  
662 GATK. Allele frequencies were calculated by dividing the number of reads supporting the alternative  
663 nucleotide with the total number of reads at each variant site.

664

#### 665 **Transcriptome evaluation**

666 Transcriptome statistics such as GC content, contig length, and assembled bases were calculated using  
667 Bash scripts and in the R statistical environment (v. 3.2.5) (Pages *et al.* 2015; R Core Team 2015). P-

668 values were corrected for multiple testing with the Benjamini and Hochberg false discovery rate  
669 (Benjamini & Hochberg 1995) and corrected values have been labelled as q-values throughout. We  
670 calculated the GC content of two *Eimeria* transcriptomes downloaded from ToxoDB (v. 25) (Gajria *et*  
671 *al.* 2007), initially sequenced by Reid *et al.* (2014). Transcriptome E90N50 was calculated using  
672 RSEM (v. 1.2.21) and Trinity (v. 2.0.6) (Li & Dewey 2011; Grabherr *et al.* 2011; Haas 2016). Plots  
673 were made with the R package ggplot2 (v. 2.2.1) (Wickham 2009). The transcriptome of *Plasmodium*  
674 *falciparum* 3D7 (v. 25) was downloaded from PlasmoDB (Gardner *et al.* 2002; Aurecochea *et al.*  
675 2009) and gene ontology information was derived from UniProtKB (Bateman *et al.* 2015). The red  
676 blood cell invasion genes were searched for in the transcriptome annotation we produced for *P.*  
677 *ashfordi* (Table S1, Supporting information). Only genes with documented involvement in  
678 *Plasmodium* erythrocyte invasion were included in the search.

679

680

## 681 **Data accessibility**

682

683 The supplementary tables and figures supporting this article have been uploaded as part of the online  
684 supporting information. The sequence reads of both host and parasite have been deposited at the  
685 NCBI Sequence Read Archive (SRA) under the accession number PRJNA311546. The assembled *P.*  
686 *ashfordi* transcriptome is available for download as Data S1 and Data S2 (Supporting information), or  
687 alternatively at <http://mbio-serv2.mbioekol.lu.se/Malavi/Downloads>.

688

689

## 690 **Author contributions**

691

692 The study design was initially conceived by OH, VP and GV, and further developed together with EV  
693 and CKC. The infection experiment was planned by GV and VP, and performed by VP. OH  
694 performed the RNA extractions. Assembly and all bioinformatic and statistical analyses were  
695 performed by EV. DA advised in the trimming of the contigs and in the mapping of sequence reads.  
696 OH, CKC, and EV planned the paper. EV wrote the paper with extensive input from all authors.

697

698

## 699 **Funding**

700

701 This work was supported by the Swedish Research Council (grant 621-2011-3548 to OH and 2010-  
702 5641 to CKC), the Crafoord Foundation (grant 20120630 to OH), European Social Fund under the  
703 global grant measure (grant VPI-3.1.-ŠMM-07-K-01-047 to GV), Research Council of Lithuania  
704 (grant MIP-045/2015 to GV) and a Wallenberg Academy Fellowship to CKC.

705

706

## 707 **Acknowledgements**

708

709 We would like to thank Staffan Bensch for stimulating discussions and comments on this paper. We  
710 are also grateful to Brian Haas for advice on transcriptome assemblies, and the director of the  
711 Biological Station “Rybachy”, Casimir V. Bolshakov, for generously providing facilities for the  
712 experimental research. Two anonymous reviewers provided valuable feedback which significantly  
713 improved the paper. The assembly and blastn computations were performed on resources provided by  
714 SNIC through Uppsala Multidisciplinary Center for Advanced Computational Science (UPPMAX)

715 (Lampa *et al.* 2013) under project b2014134; all other computational analyses were performed by EV

716 on a local machine.

717

718

## 719 **References**

720

- 721 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool.  
722 *Journal of Molecular Biology*, **215**, 403–410.
- 723 Anders S, Huber W (2010) Differential expression analysis for sequence count data. *Genome Biology*,  
724 **11**, R106.
- 725 Aurrecochea C, Brestelli J, Brunk BP *et al.* (2009) PlasmoDB: a functional genomic database for  
726 malaria parasites. *Nucleic Acids Research*, **37**, D539–D543.
- 727 Baker RP, Wijetilaka R, Urban S (2006) Two Plasmodium Rhomboid Proteases Preferentially Cleave  
728 Different Adhesins Implicated in All Invasive Stages of Malaria. *PLOS Pathogens*, **2**, e113.
- 729 Bateman A, Martin MJ, O'Donovan C *et al.* (2015) UniProt: a hub for protein information. *Nucleic*  
730 *Acids Research*, **43**, D204–D212.
- 731 Baum J, Richard D, Healer J *et al.* (2006) A conserved molecular motor drives cell invasion and  
732 gliding motility across malaria life cycle stages and other apicomplexan parasites. *Journal of*  
733 *Biological Chemistry*, **281**, 5197–5208.
- 734 Beeson JG, Drew DR, Boyle MJ *et al.* (2016) Merozoite surface proteins in red blood cell invasion,  
735 immunity and vaccines against malaria. *FEMS Microbiology Reviews*, **40**, 343–372.
- 736 Benjamini Y, Hochberg Y (1995) Controlling the False Discovery Rate: a Practical and Powerful  
737 Approach to Multiple Testing. *Journal of the Royal Statistical Society. Series B*, **57**, 289–300.
- 738 Bensch S, Canbäck B, DeBarry JD *et al.* (2016) The Genome of *Haemoproteus tartakovskyi* and Its  
739 Relationship to Human Malaria Parasites. *Genome Biology and Evolution*, **8**, 1361–1373.
- 740 Bensch S, Coltman DW, Davis CS *et al.* (2014) Genomic Resources Notes accepted 1 June 2013–31  
741 July 2013. *Molecular Ecology Resources*, **14**, 218.
- 742 Bensch S, Hellgren O, Pérez-Tris J (2009) MalAvi: a public database of malaria parasites and related  
743 haemosporidians in avian hosts based on mitochondrial cytochrome b lineages. *Molecular*  
744 *Ecology Resources*, **9**, 1353–1358.
- 745 Bensch S, Pérez-Tris J, Waldenström J, Hellgren O (2004) Linkage between nuclear and  
746 mitochondrial DNA sequences in avian malaria parasites: multiple cases of cryptic speciation?  
747 *Evolution*, **58**, 1617–1621.
- 748 Bozdech Z, Llinás M, Pulliam BL *et al.* (2003) The Transcriptome of the Intraerythrocytic  
749 Developmental Cycle of *Plasmodium falciparum*. *PLOS Biology*, **1**, e5.
- 750 Buchfink B, Xie C, Huson DH (2014) Fast and sensitive protein alignment using DIAMOND. *Nature*  
751 *Methods*, **12**, 59–60.
- 752 Camacho C, Coulouris G, Avagyan V *et al.* (2009) BLAST+: architecture and applications. *BMC*  
753 *Bioinformatics*, **10**, 421.
- 754 Cingolani P, Platts A, Wang LL *et al.* (2012) A program for annotating and predicting the effects of  
755 single nucleotide polymorphisms, SnpEff. *Fly*, **6**, 80–92.
- 756 Cordery DV, Kishore U, Kyes S *et al.* (2007) Characterization of a *Plasmodium falciparum*  
757 Macrophage-Migration Inhibitory Factor Homologue. *The Journal of Infectious Diseases*, **195**,  
758 905–912.
- 759 Cornet S, Bichet C, Larcombe S, Faivre B, Sorci G (2014) Impact of host nutritional status on  
760 infection dynamics and parasite virulence in a bird-malaria system. *Journal of Animal Ecology*,  
761 **83**, 256–265.
- 762 Counihan NA, Kalanon M, Coppel RL, De Koning-Ward TF (2013) *Plasmodium* rhoptry proteins:  
763 why order is important. *Trends in Parasitology*, **29**, 228–236.

- 764 Cowman AF, Berry D, Baum J (2012) The cellular and molecular basis for malaria parasite invasion  
765 of the human red blood cell. *Journal of Cell Biology*, **198**, 961–971.
- 766 Craig A, Scherf A (2001) Molecules on the surface of the Plasmodium falciparum infected  
767 erythrocyte and their role in malaria pathogenesis and immune evasion. *Molecular and*  
768 *Biochemical Parasitology*, **115**, 129–143.
- 769 Daily JP, Le Roch KG, Sarr O *et al.* (2005) In vivo transcriptome of Plasmodium falciparum reveals  
770 overexpression of transcripts that encode surface proteins. *Journal of Infectious Diseases*, **191**,  
771 1196–1203.
- 772 Daily JP, Scanfeld D, Pochet N *et al.* (2007) Distinct physiological states of Plasmodium falciparum  
773 in malaria-infected patients. *Nature*, **450**, 1091–1095.
- 774 Dimitrov D, Palinauskas V, Iezhova TA *et al.* (2015) Plasmodium spp.: An experimental study on  
775 vertebrate host susceptibility to avian malaria. *Experimental Parasitology*, **148**, 1–16.
- 776 Drovetski S V, Aghayan SA, Mata VA *et al.* (2014) Does the niche breadth or trade-off hypothesis  
777 explain the abundance-occupancy relationship in avian Haemosporidia? *Molecular Ecology*, **23**,  
778 3322–3329.
- 779 Ellis VA, Cornet S, Merrill L *et al.* (2015) Host immune responses to experimental infection of  
780 Plasmodium relictum (lineage SGS1) in domestic canaries (Serinus canaria). *Parasitology*  
781 *Research*, **114**, 3627–3636.
- 782 Farias ME, Atkinson CT, LaPointe DA, Jarvi SI (2012) Analysis of the trap gene provides evidence  
783 for the role of elevation and vector abundance in the genetic diversity of Plasmodium relictum in  
784 Hawaii. *Malaria Journal*, **11**, 305.
- 785 Gajria B, Bahl A, Brestelli J *et al.* (2007) ToxoDB: An integrated toxoplasma gondii database  
786 resource. *Nucleic Acids Research*, **36**, D553–D556.
- 787 Gardner MJ, Hall N, Fung E *et al.* (2002) Genome sequence of the human malaria parasite  
788 Plasmodium falciparum. *Nature*, **419**, 498–511.
- 789 Garnham PCC (1966) *Malaria parasites and other Haemosporidia*. Blackwell Scientific Publications  
790 Ltd., Oxford, UK.
- 791 Grabherr MG, Haas BJ, Yassour M *et al.* (2011) Full-length transcriptome assembly from RNA-Seq  
792 data without a reference genome. *Nature Biotechnology*, **29**, 644–652.
- 793 Grech K, Watt K, Read AF (2006) Host-parasite interactions for virulence and resistance in a malaria  
794 model system. *Journal of Evolutionary Biology*, **19**, 1620–1630.
- 795 Green JL, Hinds L, Grainger M, Knuepfer E, Holder AA (2006) Plasmodium thrombospondin related  
796 apical merozoite protein (PTRAMP) is shed from the surface of merozoites by PfSUB2 upon  
797 invasion of erythrocytes. *Molecular and Biochemical Parasitology*, **150**, 114–117.
- 798 Guizetti J, Scherf A (2013) Silence, activate, poise and switch! Mechanisms of antigenic variation in  
799 Plasmodium falciparum. *Cellular Microbiology*, **15**, 718–726.
- 800 Haas B (2016) Transcriptome Contig Nx and ExN50 stats.  
801 <https://github.com/trinityrnaseq/trinityrnaseq/wiki/Transcriptome-Contig-Nx-and-ExN50-stats>.
- 802 Hall N, Karras M, Raine JD *et al.* (2005) A Comprehensive Survey of the Plasmodium Life Cycle by  
803 Genomic, Transcriptomic, and Proteomic Analyses. *Science*, **307**, 82–86.
- 804 Hellgren O, Atkinson CT, Bensch S *et al.* (2015) Global phylogeography of the avian malaria  
805 pathogen Plasmodium relictum based on MSP1 allelic diversity. *Ecography*, **38**, 842–850.
- 806 Hellgren O, Kutzer M, Bensch S, Valkiūnas G, Palinauskas V (2013) Identification and  
807 characterization of the merozoite surface protein 1 (msp1) gene in a host-generalist avian  
808 malaria parasite, Plasmodium relictum (lineages SGS1 and GRW4) with the use of blood  
809 transcriptome. *Malaria Journal*, **12**, 381.
- 810 Hellgren O, Waldenström J, Bensch S (2004) A new PCR assay for simultaneous studies of

- 811 Leucocytozoon, Plasmodium, and Haemoproteus from avian blood. *Journal of Parasitology*, **90**,  
812 797–802.
- 813 Hisaeda H, Maekawa Y, Iwakawa D *et al.* (2004) Escape of malaria parasites from host immunity  
814 requires CD4+CD25+ regulatory T cells. *Nature Medicine*, **10**, 29–30.
- 815 Idaghdour Y, Quinlan J, Goulet J-P *et al.* (2012) Evidence for additive and interaction effects of host  
816 genotype and infection in malaria. *Proceedings of the National Academy of Sciences of the*  
817 *United States of America*, **109**, 16786–16793.
- 818 Jeffares DC, Pain A, Berry A *et al.* (2007) Genome variation and evolution of the malaria parasite  
819 *Plasmodium falciparum*. *Nature Genetics*, **39**, 120–125.
- 820 Josling GA, Petter M, Oehring SC *et al.* (2015) A Plasmodium Falciparum Bromodomain Protein  
821 Regulates Invasion Gene Expression. *Cell Host & Microbe*, **17**, 741–751.
- 822 Kersey PJ, Allen JE, Armean I *et al.* (2016) Ensembl Genomes 2016: More genomes, more  
823 complexity. *Nucleic Acids Research*, **44**, D574–D580.
- 824 Koutsovoulos G, Kumar S, Laetsch DR *et al.* (2016) No evidence for extensive horizontal gene  
825 transfer in the genome of the tardigrade *Hypsibius dujardini*. *Proceedings of the National*  
826 *Academy of Sciences of the United States of America*, **113**, 5053–5058.
- 827 Križanauskienė A, Hellgren O, Kosarev V *et al.* (2006) Variation in host specificity between species  
828 of avian hemosporidian parasites: evidence from parasite morphology and cytochrome B gene  
829 sequences. *Journal of Parasitology*, **92**, 1319–1324.
- 830 Lachish S, Knowles SCL, Alves R, Wood MJ, Sheldon BC (2011) Fitness effects of endemic malaria  
831 infections in a wild bird population: The importance of ecological structure. *Journal of Animal*  
832 *Ecology*, **80**, 1196–1206.
- 833 Lampa S, Dahlö M, Olason P, Hagberg J, Spjuth O (2013) Lessons learned from implementing a  
834 national infrastructure in Sweden for storage and analysis of next-generation sequencing data.  
835 *GigaScience*, **2**, 9.
- 836 Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. *Nature Methods*, **9**,  
837 357–359.
- 838 Lapp SA, Mok S, Zhu L *et al.* (2015) Plasmodium knowlesi gene expression differs in ex vivo  
839 compared to in vitro blood-stage cultures. *Malaria Journal*, **14**, 110.
- 840 Lauron EJ, Aw Yeang HX, Taffner SM, Sehgal RNM (2015) De novo assembly and transcriptome  
841 analysis of Plasmodium gallinaceum identifies the Rh5 interacting protein (ripr), and reveals a  
842 lack of EBL and RH gene family diversification. *Malaria Journal*, **14**, 296.
- 843 Lauron EJ, Oakgrove KS, Tell LA *et al.* (2014) Transcriptome sequencing and analysis of  
844 Plasmodium gallinaceum reveals polymorphisms and selection on the apical membrane antigen-  
845 1. *Malaria Journal*, **13**, 382.
- 846 Lemieux JE, Gomez-Escobar N, Feller A *et al.* (2009) Statistical estimation of cell-cycle progression  
847 and lineage commitment in Plasmodium falciparum reveals a homogeneous pattern of  
848 transcription in ex vivo culture. *Proceedings of the National Academy of Sciences of the United*  
849 *States of America*, **106**, 7559–7564.
- 850 LeRoux M, Lakshmanan V, Daily JP (2009) Plasmodium falciparum biology: analysis of in vitro  
851 versus in vivo growth conditions. *Trends in Parasitology*, **25**, 474–481.
- 852 Li B, Dewey CN (2011) RSEM: accurate transcript quantification from RNA-Seq data with or  
853 without a reference genome. *BMC Bioinformatics*, **12**, 323.
- 854 Li W, Godzik A (2006) Cd-hit: A fast program for clustering and comparing large sets of protein or  
855 nucleotide sequences. *Bioinformatics*, **22**, 1658–1659.
- 856 Love MI, Huber W, Anders S (2014) Moderated estimation of fold change and dispersion for RNA-  
857 seq data with DESeq2. *Genome Biology*, **15**, 550.



- 858 Lovegrove FE, Peña-Castillo L, Mohammad N *et al.* (2006) Simultaneous host and parasite  
859 expression profiling identifies tissue-specific transcriptional programs associated with  
860 susceptibility or resistance to experimental cerebral malaria. *BMC Genomics*, **7**, 295.
- 861 Lutz HL, Marra NJ, Grewe F *et al.* (2016) Laser capture microdissection microscopy and genome  
862 sequencing of the avian malaria parasite, *Plasmodium relictum*. *Parasitology Research*, **115**,  
863 4503–4510.
- 864 Mackinnon M, Gaffney D, Read A (2002) Virulence in rodent malaria: host genotype by parasite  
865 genotype interactions. *Infection, Genetics and Evolution*, **1**, 287–296.
- 866 Martinez C, Marzec T, Smith CD, Tell LA, Sehgal RNM (2013) Identification and expression of  
867 *maebl*, an erythrocyte-binding gene, in *Plasmodium gallinaceum*. *Parasitology Research*, **112**,  
868 945–954.
- 869 Martinsen ES, Perkins SL (2013) The diversity of *Plasmodium* and other haemosporidians: The  
870 intersection of taxonomy, phylogenetics and genomics. In: *Malaria parasites: comparative*  
871 *genomics, evolution and molecular biology*, pp. 1–15. Caister Academic Press, Norfolk.
- 872 McKenna A, Hanna M, Banks E *et al.* (2010) The Genome Analysis Toolkit: A MapReduce  
873 framework for analyzing next-generation DNA sequencing data. *Genome Research*, **20**, 1297–  
874 1303.
- 875 Mueller C, Klages N, Jacot D *et al.* (2013) The *Toxoplasma* Protein ARO Mediates the Apical  
876 Positioning of Rhoptry Organelles, a Prerequisite for Host Cell Invasion. *Cell Host & Microbe*,  
877 **13**, 289–301.
- 878 Niang M, Bei AK, Madnani KG *et al.* (2014) STEVOR Is a *Plasmodium falciparum* Erythrocyte  
879 Binding Protein that Mediates Merozoite Invasion and Rosetting. *Cell Host & Microbe*, **16**, 81–  
880 93.
- 881 Opitz C, Soldati D (2002) “The glideosome”: a dynamic complex powering gliding motion and host  
882 cell invasion by *Toxoplasma gondii*. *Molecular Microbiology*, **45**, 597–604.
- 883 Otto TD, Böhme U, Jackson AP *et al.* (2014a) A comprehensive evaluation of rodent malaria parasite  
884 genomes and gene expression. *BMC Biology*, **12**, 86.
- 885 Otto TD, Rayner JC, Böhme U *et al.* (2014b) Genome sequencing of chimpanzee malaria parasites  
886 reveals possible pathways of adaptation to human hosts. *Nature Communications*, **5**, 4754.
- 887 Otto TD, Wilinski D, Assefa S *et al.* (2010) New insights into the blood-stage transcriptome of  
888 *Plasmodium falciparum* using RNA-Seq. *Molecular Microbiology*, **76**, 12–24.
- 889 Pages H, Aboyoun P, Gentleman R, DebRoy S (2015) Biostrings: String objects representing  
890 biological sequences, and matching algorithms. *R package version 2.38.4*.
- 891 Palinauskas V, Valkiūnas G, Bolshakov CV, Bensch S (2008) *Plasmodium relictum* (lineage P-  
892 SGS1): effects on experimentally infected passerine birds. *Experimental Parasitology*, **120**,  
893 372–380.
- 894 Palinauskas V, Valkiūnas G, Bolshakov CV, Bensch S (2011) *Plasmodium relictum* (lineage SGS1)  
895 and *Plasmodium ashfordi* (lineage GRW2): the effects of the co-infection on experimentally  
896 infected passerine birds. *Experimental Parasitology*, **127**, 527–533.
- 897 Pauli M, Chakarov N, Rupp O *et al.* (2015) De novo assembly of the dual transcriptomes of a  
898 polymorphic raptor species and its malarial parasite. *BMC Genomics*, **16**, 1038.
- 899 Pérez-Tris J, Hellgren O, Križanauskienė A *et al.* (2007) Within-host speciation of malaria parasites.  
900 *PLOS ONE*, **2**, e235.
- 901 R Core Team (2015) R: A language and environment for statistical computing. *R Foundation for*  
902 *Statistical Computing, Vienna, Austria*.
- 903 Recker M, Buckee CO, Serazin A *et al.* (2011) Antigenic Variation in *Plasmodium falciparum*  
904 Malaria Involves a Highly Structured Switching Pattern. *PLOS Pathogens*, **7**, e1001306.

- 905 Reid AJ, Blake DP, Ansari HR *et al.* (2014) Genomic analysis of the causative agents of coccidiosis  
906 in domestic chickens. *Genome Research*, **24**, 1676–1685.
- 907 de Roode JC, Culleton R, Cheesman SJ, Carter R, Read AF (2004) Host heterogeneity is a  
908 determinant of competitive exclusion or coexistence in genetically diverse malaria infections.  
909 *Proceedings of the Royal Society B: Biological Sciences*, **271**, 1073–1080.
- 910 Santos JM, Graindorge A, Soldati-Favre D (2012) New insights into parasite rhomboid proteases.  
911 *Molecular and Biochemical Parasitology*, **182**, 27–36.
- 912 Schmieder R, Edwards R (2011) Quality control and preprocessing of metagenomic datasets.  
913 *Bioinformatics*, **27**, 863–864.
- 914 Siau A, Silvie O, Franetich J-F *et al.* (2008) Temperature Shift and Host Cell Contact Up-Regulate  
915 Sporozoite Expression of Plasmodium falciparum Genes Involved in Hepatocyte Infection.  
916 *PLOS Pathogens*, **4**, e1000121.
- 917 Siegel T, Hon C-C, Zhang Q *et al.* (2014) Strand-specific RNA-Seq reveals widespread and  
918 developmentally regulated transcription of natural antisense transcripts in Plasmodium  
919 falciparum. *BMC Genomics*, **15**, 150.
- 920 Silamut K, Phu NH, Whitty C *et al.* (1999) A Quantitative Analysis of the Microvascular  
921 Sequestration of Malaria Parasites in the Human Brain. *American Journal of Pathology*, **155**,  
922 395–410.
- 923 Spence PJ, Jarra W, Lévy P *et al.* (2013) Vector transmission regulates immune control of  
924 Plasmodium virulence. *Nature*, **498**, 228–231.
- 925 Templeton TJ, Kaslow DC (1997) Cloning and cross-species comparison of the thrombospondin-  
926 related anonymous protein (TRAP) gene from Plasmodium knowlesi, Plasmodium vivax and  
927 Plasmodium gallinaceum. *Molecular and Biochemical Parasitology*, **84**, 13–24.
- 928 Valkiūnas G, Zehtindjiev P, Hellgren O *et al.* (2007) Linkage between mitochondrial cytochrome b  
929 lineages and morphospecies of two avian malaria parasites, with a description of Plasmodium  
930 (Novyella) ashfordi sp. nov. *Parasitology Research*, **100**, 1311–1322.
- 931 Videvall E, Cornwallis CK, Palinauskas V, Valkiūnas G, Hellgren O (2015) The Avian Transcriptome  
932 Response to Malaria Infection. *Molecular Biology and Evolution*, **32**, 1255–1267.
- 933 Vincensini L, Fall G, Berry L, Blisnick T, Braun Breton C (2008) The RhopH complex is transferred  
934 to the host cell cytoplasm following red blood cell invasion by Plasmodium falciparum.  
935 *Molecular and Biochemical Parasitology*, **160**, 81–89.
- 936 Wickham H (2009) *ggplot2: elegant graphics for data analysis*. New York: Springer.
- 937 Young JA, Fivelman QL, Blair PL *et al.* (2005) The Plasmodium falciparum sexual development  
938 transcriptome: A microarray analysis using ontology-based pattern identification. *Molecular and*  
939 *Biochemical Parasitology*, **143**, 67–79.
- 940 Zehtindjiev P, Ilieva M, Westerdahl H *et al.* (2008) Dynamics of parasitemia of malaria parasites in a  
941 naturally and experimentally infected migratory songbird, the great reed warbler *Acrocephalus*  
942 *arundinaceus*. *Experimental Parasitology*, **119**, 99–110.
- 943

944 **Tables**

945

946

947 **Table 1.** Assembly statistics of the *Plasmodium ashfordi* transcriptome.

948

	Annotated transcripts	Unannotated transcripts	Total assembly
Number of contigs	7 860	4 094	11 954
Number of bases	7 316 007	1 694 373	9 010 380
Contig length min-max (bp)	200 – 26 773	200 – 4 171	200 – 26 773
Median contig length (bp)	681.0	321.0	498.0
Mean contig length (bp)	930.8	413.9	753.8
GC content (%)	21.22	17.26	20.48

949

950

951

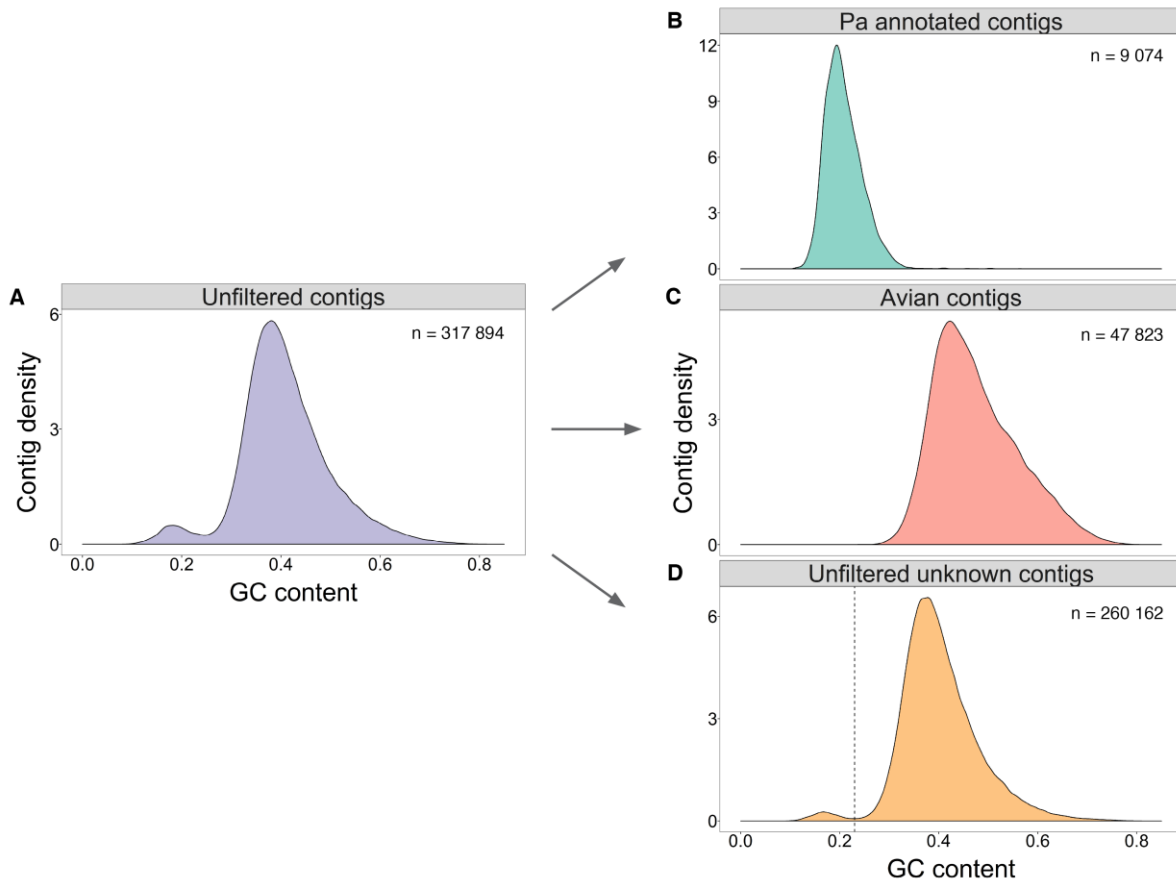
952 **Table 2.** Assembled transcripts of genes involved in *Plasmodium* invasion of red blood cells.

953

Gene name	Gene product	Represented by transcripts	Species match (blastx)	Bit score	e-value
ACT1	Actin 1	TR215622 c0_g1_i1	<i>P. vivax</i>	758.4	5.5E-216
ALP1	Actin-like protein 1	TR55613 c0_g1_i1	<i>P. vivax</i>	619	5.1E-174
AMA-1	Apical membrane antigen 1	TR118224 c0_g1_i1	<i>P. gallinaceum</i>	624	2.1E-175
ARO	Armadillo-domain containing rhopty protein	TR18124 c0_g1_i1 TR18124 c0_g2_i1	<i>P. knowlesi</i>	520 501.9	1.9E-144 6.0E-139
BDP1	Bromodomain protein 1	TR66841 c0_g2_i1 TR66841 c0_g2_i2	<i>P. falciparum</i>	364.8 357.1	1.9E-97 3.9E-95
FBPA	Fructose-bisphosphate aldolase	TR12911 c0_g1_i1 TR12911 c0_g2_i1	<i>P. cynomolgi</i> <i>P. berghei</i>	223.8 353.2	2.5E-55 2.7E-94
GAMA	GPI-anchored micronemal antigen	TR16322 c0_g2_i1	<i>P. reichenowi</i>	567.4	2.5E-158
GAP45	Glideosome-associated protein 45	TR34884 c0_g1_i1	<i>P. vivax</i>	209.9	3.9E-51
GAP50	Glideosome-associated protein 50	TR144721 c0_g1_i1	<i>P. knowlesi</i>	577.8	1.2E-161
MAEBL	Merozoite adhesive erythrocytic binding protein	TR234951 c0_g1_i1 TR99718 c0_g1_i1 TR125315 c0_g1_i1	<i>P. yoelii</i> <i>P. gallinaceum</i> <i>P. gallinaceum</i>	171.8 146.4 84.3	6.3E-40 2.4E-32 7.2E-14
MCP1	Merozoite capping protein 1	TR122100 c0_g1_i1	<i>P. berghei</i>	164.5	1.7E-37
MSP1	Merozoite surface protein 1	TR112241 c0_g2_i1 TR176579 c0_g1_i1	<i>P. relictum</i>	368.6 134.4	4.2E-98 1.4E-28
MTIP (MLC1)	Myosin A tail domain interacting protein	TR8937 c0_g1_i1	<i>P. vivax</i>	295	8.6E-77
MTRAP	Merozoite TRAP-like protein	TR188691 c0_g2_i1	<i>P. vinckei</i>	98.6	1.3E-17
MYOA	Myosin A	TR198550 c0_g1_i1	<i>Babesia microti</i>	199.1	1.5E-47
RAMA	Rhoptry-associated membrane antigen	TR144799 c0_g1_i1	<i>P. knowlesi</i>	70.1	5.0E-09
RALP1	Rhoptry-associated leucine zipper-like protein 1	TR7446 c0_g1_i1	<i>P. falciparum</i>	102.1	1.2E-18
RAP1	Rhoptry-associated protein 1	TR115690 c2_g1_i1	<i>P. coatneyi</i>	236.1	1.3E-58
RAP3	Rhoptry-associated protein 3	TR13305 c0_g1_i1	<i>P. falciparum</i>	202.6	9.9E-49
RBP1 (RH1)	Reticulocyte-binding protein 1	TR53756 c1_g1_i1	<i>P. falciparum</i>	84.7	8.0E-13
RBP2 (RH2)	Reticulocyte-binding protein 2	TR208311 c0_g1_i1 TR66282 c1_g2_i1	<i>P. vivax</i>	76.6 146.7	5.9E-11 2.1E-31
RBP2b (RH2b)	Reticulocyte-binding protein 2 homolog B	TR87235 c1_g1_i1 TR35437 c1_g1_i1	<i>P. vivax</i> <i>P. falciparum</i>	97.8 86.7	8.7E-18 1.9E-14
RHOPH1	High molecular weight rhopty protein 1 (CLAG9)	TR7367 c0_g1_i1 TR200056 c0_g1_i1 TR233854 c0_g1_i1	<i>P. reichenowi</i> <i>P. reichenowi</i> <i>P. falciparum</i>	272.3 82.8 61.6	7.9E-70 2.3E-13 4.8E-07
RHOPH2	High molecular weight rhopty protein 2	TR25858 c0_g2_i1 TR25858 c0_g3_i1 TR175810 c0_g1_i1	<i>P. vivax</i> <i>P. fragile</i> <i>P. vivax</i>	542.3 335.9 180.6	7.8E-151 5.7E-89 3.4E-42
RHOPH3	High molecular weight rhopty protein 3	TR83052 c6_g1_i1	<i>P. falciparum</i>	612.8	7.3E-172
RIPR	RH5 interacting protein	TR49727 c0_g2_i1 TR49727 c0_g2_i2	<i>P. reichenowi</i> <i>P. fragile</i>	502.3 299.3	6.4E-139 1.0E-77
ROM1	Rhomboid protease 1	TR178976 c0_g1_i1	<i>P. yoelii</i>	102.4	2.4E-19
ROM4	Rhomboid protease 4	TR92176 c1_g1_i1	<i>P. reichenowi</i>	783.5	2.1E-223
RON2	Rhoptry neck protein 2	TR62454 c0_g1_i1	<i>P. cynomolgi</i>	659.8	8.5E-186
RON4	Rhoptry neck protein 4	TR145809 c0_g1_i1	<i>P. reichenowi</i>	349.4	8.9E-93
RON5	Rhoptry neck protein 5	TR67313 c1_g1_i1	<i>P. reichenowi</i>	1105.9	0
SUB1	Subtilisin 1	TR164664 c0_g1_i1	<i>P. inui</i>	636.7	3.0E-179
SUB2	Subtilisin 2	TR169322 c0_g1_i1	<i>P. inui</i>	563.9	2.7E-157
TLP	TRAP-like protein	TR153597 c0_g1_i1	<i>P. falciparum</i>	440.7	3.2E-120
TRAMP	Thrombospondin related apical membrane protein	TR179194 c0_g1_i1	<i>P. falciparum</i>	231.1	2.1E-57
TRAP	Thrombospondin related anonymous protein	TR16495 c0_g1_i1	<i>P. relictum</i>	229.2	6.4E-57



## 955 Figures



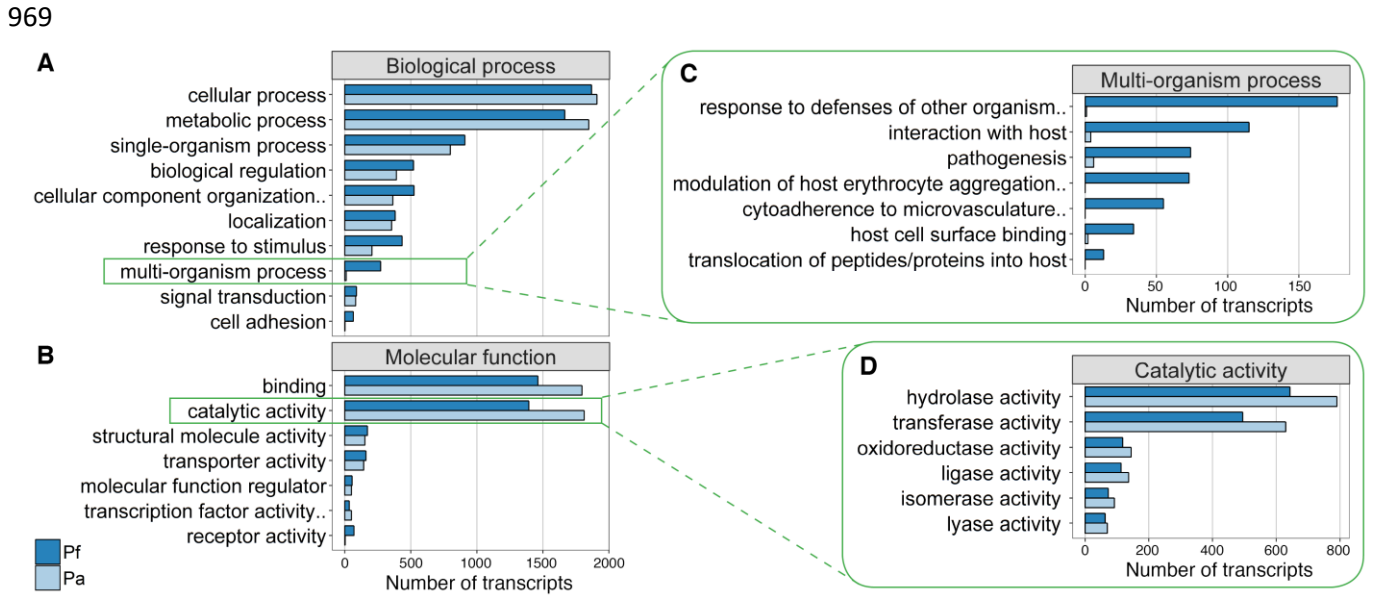
956

957

958 **Figure 1.** Filtering of host transcripts from the *Plasmodium ashfordi* transcriptome using gene  
959 annotation and GC content. Density curves of contig GC content in (A) the initial, unfiltered  
960 assembly, (B) the annotated *P. ashfordi* transcriptome assembly, (C) all contigs giving significant  
961 blastx matches to birds, and (D) all unknown contigs before GC % filtering. The arrows indicate  
962 assembly versions before and after initial filtering and cleaning steps. Both the initial, unfiltered  
963 assembly and the assembly with unknown, unfiltered contigs display a bimodal distribution,  
964 incorporating both avian and malaria parasite transcripts. The dashed straight line in D indicates the  
965 23% GC cut-off where unknown transcripts with lower GC content were extracted, filtered, and later  
966 included in the final *P. ashfordi* assembly as unannotated transcripts.

967

968



970

971

972 **Figure 2.** Gene ontology terms of transcripts in the *Plasmodium ashfordi* transcriptome assembly (Pa; light blue) compared to the transcriptome of the human parasite *P. falciparum* (Pf; dark blue).

973 Information of gene ontology was successfully retrieved for 4 346 annotated *P. ashfordi* transcripts and 3 751 *P. falciparum* transcripts. Gene ontology terms containing at least 50 transcripts in one of

974 the species under the two major categories (A) 'biological process' and (B) 'molecular function' are shown. Details of underlying child terms are given for two categories: (C) 'multi-organism process'

975 and (D) 'catalytic activity', where *P. ashfordi* displays fewer and more transcripts, respectively,

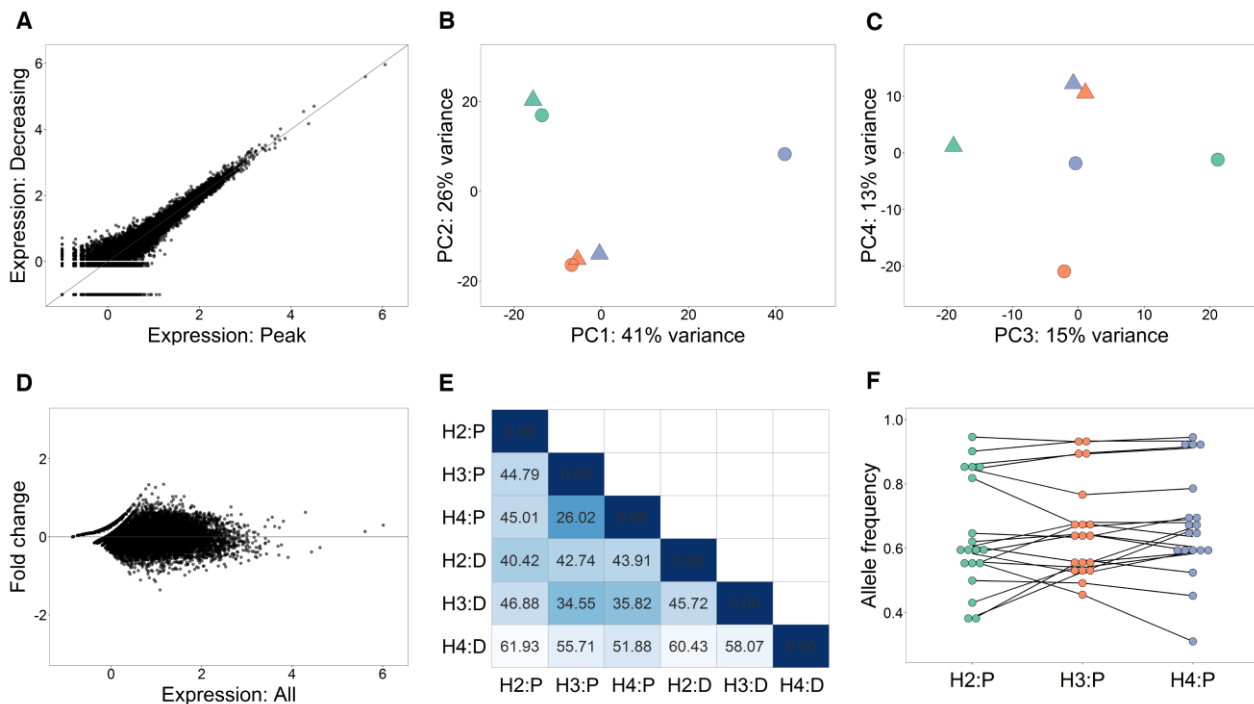
976 relative *P. falciparum*. Gene ontology terms containing at least 10 transcripts in one of the species are shown in C-D. Names of terms ending with dots have been shortened in the figure due to space

977 limitations. For a complete list of terms, see Tables S8-S14 (Supporting information).

978

982

983



984

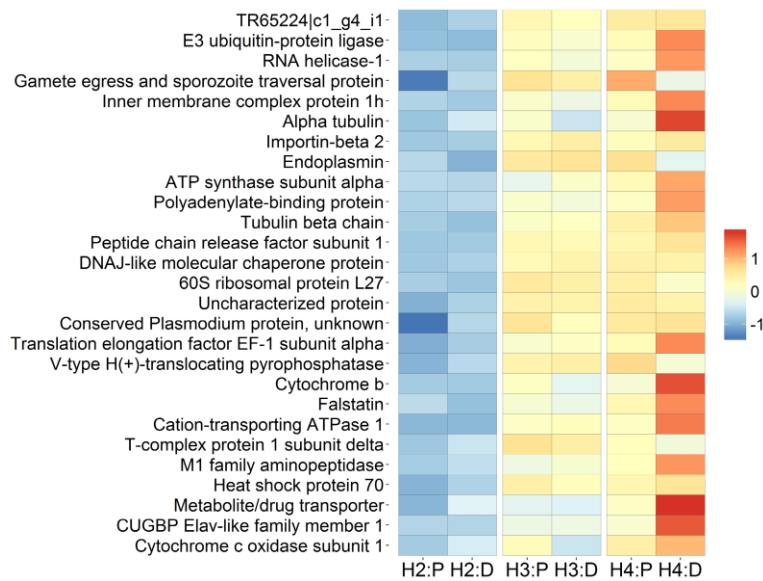
985

986 **Figure 3.** Gene expression patterns of *Plasmodium ashfordi* in individual hosts during two infection  
 987 stages. **(A)** Scatter plot displaying expression levels of all transcripts in the *P. ashfordi* transcriptome  
 988 ( $n = 11\,954$ ). The axes show log-transformed normalized mean expression values + 0.1 during peak  
 989 parasitemia stage (x-axis;  $n = 3$ ) and decreasing parasitemia stage (y-axis;  $n = 3$ ). **(B-C)** Principal  
 990 component analysis (PCA) plots show clustering of samples based on variation in regularized log-  
 991 transformed normalized gene expression. **(B)** shows principal component 1 and 2, and **(C)** shows  
 992 principal component 3 and 4. Colours of parasite transcriptomes illustrate which host they are  
 993 sampled from: host 2 = green, host 3 = orange, host 4 = purple. Triangles and circles, respectively,  
 994 indicate parasites sampled during peak and decreasing parasitemia stages. **(D)** MA plot showing log-  
 995 transformed normalized expression values + 0.1 in *P. ashfordi* averaged over all samples (x-axis;  $n =$   
 996 6) and shrunken  $\log_2$  expression fold changes between the parasitemia stages (y-axis). **(E)** Heatmap  
 997 portraying Euclidian distance measures between parasite expression profiles in different hosts during  
 998 the two time points. Lighter colour signifies greater distance. H = host, P = peak parasitemia, and D =  
 999 decreasing parasitemia. The distances between parasite transcriptomes in the decreasing parasitemia  
 1000 stage to any transcriptome sampled during peak parasitemia are shortest within the same host  
 1001 individual. **(F)** Dot plot illustrating allele frequencies of 19 SNPs in the *P. ashfordi* transcriptome  
 1002 from three hosts during peak parasitemia stage. Black lines connect the same SNP for comparisons  
 1003 between hosts. All SNPs contained identical alleles in the parasites from the different hosts and  
 1004 displayed no differences in allele frequency between the host individuals ( $p$ -value = 0.432).

1005



1006



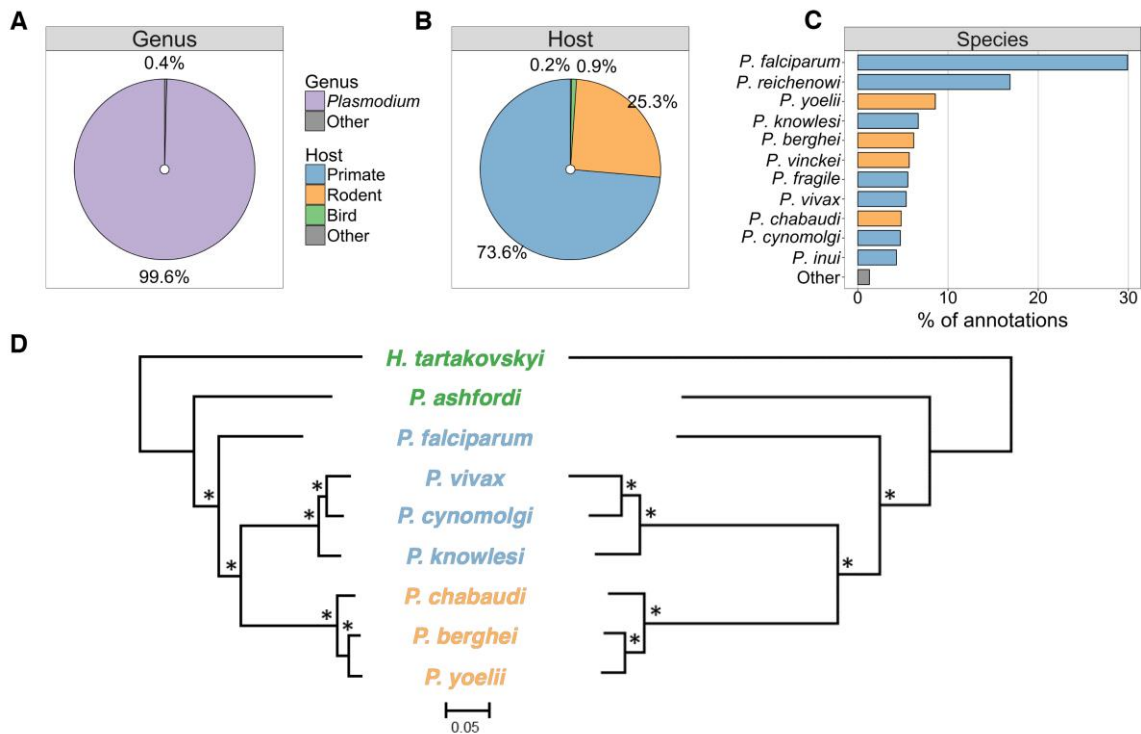
1007

1008

1009 **Figure 4.** Heatmap of relative expression levels of 28 *Plasmodium ashfordi* transcripts (rows) that  
1010 were significantly differentially expressed between parasites from different hosts (columns). Warmer  
1011 colour signifies higher transcript expression, and blue colour indicates lower expression. H = host, P =  
1012 peak parasitemia, and D = decreasing parasitemia. To compare across genes, expression levels have  
1013 been normalized with respect to library size, regularized log-transformed, and scaled and centred  
1014 around zero to give Z-scores.

1015

1016



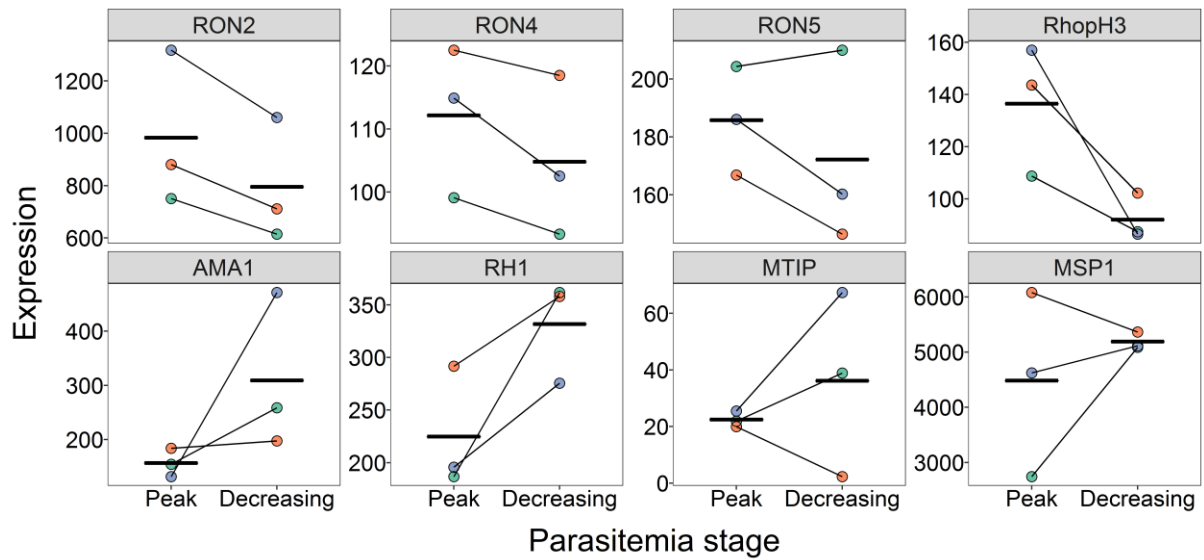
1017

1018

1019 **Figure 5.** Distribution of apicomplexan parasites presenting best sequence matches to the  
 1020 *Plasmodium ashfordi* transcriptome. **(A)** Pie chart showing the distribution of *P. ashfordi* annotations  
 1021 resulting in significant blast matches to *Plasmodium* species and parasites of other genera. **(B)**  
 1022 pie chart showing the distribution of annotations resulting in significant matches to parasites infecting  
 1023 primates, rodents, birds, and other hosts. **(C)** Bar plot displaying the proportion of annotated *P.*  
 1024 *ashfordi* contigs giving significant matches to parasites on a species level. The ‘other’ category  
 1025 contains here all other apicomplexan species, including bird parasites, comprising a total of 20  
 1026 different species (complete list can be found in Table S6, Supporting information). **(D)** Phylogenetic  
 1027 trees of *Plasmodium* parasites, rooted with the related avian blood parasite *Haemoproteus*  
 1028 *tartakovskyi*. The left tree is constructed from 599 concatenated genes shared with other species in the  
 1029 Apicomplexa phylum. The right tree is based on 703 concatenated genes that are unique to parasites  
 1030 within these two genera. \* = 100% bootstrap support. The phylogenies are derived and adapted with  
 1031 permission from Bensch *et al.* (2016). Same colour scheme as in B applies for C and D, i.e. blue  
 1032 signifies a parasite of primates, orange a parasite of rodents, and green a parasite of birds.

1033

1034



1035

1036

1037

1038

1039

1040

1041

**Figure 6.** Individual gene expression plots for some of the *Plasmodium ashfordi* transcripts involved in red blood cell invasion. Line plots displaying normalized parasite gene expression in each individual host over the two sampled parasitemia stages. Host 2 is depicted in green, host 3 in orange, and host 4 in purple. Thick horizontal lines indicate mean expression levels in each stage.

## 1042 **Supporting Information**

1043

1044

1045 **Figure S1.** Length distribution of contigs in the *Plasmodium ashfordi* transcriptome assembly.

1046

1047 **Figure S2.** Multidensity plot showing density of transcripts over log-transformed normalized  
1048 expression values.

1049

1050 **Figure S3.** Dot plot showing allele frequencies of 19 SNPs in the parasites of three hosts during peak  
1051 parasitemia.

1052

1053 **Figure S4.** GC content distribution of a transcriptome assembly of the chicken parasite *Plasmodium*  
1054 *gallinaceum* (Lauron *et al.* 2015) highlights the difficulties in assembling clean parasite  
1055 transcriptomes from dual RNA-seq data.

1056

1057 **Figure S5.** Example of transcript sequence homogeneity from the *P. ashfordi* transcript  
1058 TR213315|c0\_g1\_i1 derived from the mitochondrial gene cytochrome c oxidase subunit 1 (COX1).

1059

1060 **Figure S6.** Example of transcript sequence heterogeneity showing the presence of three SNPs in the  
1061 *P. ashfordi* transcript TR73260|c0\_g1\_i1.

1062

1063

1064 **Table S1.** Information of all annotated *P. ashfordi* transcripts (n = 7 860).

1065

1066 **Table S2.** Normalized expression levels of all *P. ashfordi* transcripts (n = 11 954) in individual hosts  
1067 during peak and decreasing parasitemia stages.

1068

1069 **Table S3.** *P. ashfordi* transcripts that were significantly differentially expressed between host  
1070 individuals (n = 28).

1071

1072 **Table S4.** Results of SNPs discovered in the *P. ashfordi* transcriptome (n = 19).

1073

1074 **Table S5.** Gametocyte and meront proportions in the *P. ashfordi* blood smears.

1075

1076 **Table S6.** Species distribution matches from the annotated transcripts of *P. ashfordi*.

1077

1078 **Table S7.** Most highly expressed transcripts in the *P. ashfordi* transcriptome.

1079

1080 **Table S8.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1081 in the category 'Biological process'.

1082

1083 **Table S9.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1084 in the category ‘Molecular function’.

1085

1086 **Table S10.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1087 in the category ‘Multi-organism process’.

1088

1089 **Table S11.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1090 in the category ‘Catalytic activity’.

1091

1092 **Table S12.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1093 in the category ‘Metabolic process’.

1094

1095 **Table S13.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1096 in the category ‘Cellular process’.

1097

1098 **Table S14.** Gene ontology terms with associated numbers of *P. ashfordi* and *P. falciparum* transcripts  
1099 in the category ‘Binding’.

1100

1101

1102 **Data S1.** Sequences in the annotated *P. ashfordi* transcriptome assembly (n = 7 860) as a gzipped  
1103 multi-fasta file.

1104

1105 **Data S2.** Sequences in the total *P. ashfordi* transcriptome assembly (n = 11 954) as a gzipped multi-  
1106 fasta file.

1107