

1 Probabilistic data integration identifies reliable gametocyte-
2 specific proteins and transcripts in malaria parasites

3
4
5
6
7
8
9
10
11

Lisette Meerstein-Kessel^{1,2}, Robin van der Lee^{1#a}, Will Stone^{2,3}, Kjerstin Lanke², David A Baker⁴, Pietro Alano⁵, Francesco Silvestrini⁵, Chris J Janse⁶, Shahid M Khan⁶, Marga van de Vegte-Bolmer², Wouter Graumans², Rianne Siebelink-Stoter², Taco WA Kooij², Matthias Marti⁷, Chris Drakeley³, Joseph J. Campo⁸, Teunis JP van Dam^{1#b}, Robert Sauerwein², Teun Bousema^{2¶}, Martijn A Huynen^{1¶*}

¹ Centre for Molecular and Biomolecular Informatics, Radboud Institute for Molecular Life Sciences, Radboud university medical center, Nijmegen, The Netherlands

² Department of Medical Microbiology, Radboud university medical center, Nijmegen, The Netherlands

³ Department of Immunology and Infection, Faculty of Infectious and Tropical Diseases, London School of Hygiene & Tropical Medicine, London, United Kingdom

⁴ Department of Pathogen Molecular Biology, Faculty of Infectious and Tropical Diseases, London School of Hygiene & Tropical Medicine, London, United Kingdom

⁵ Dipartimento Malattie Infettive, Istituto Superiore di Sanità, Rome, Italy

⁶ Department of Parasitology, Leiden University Medical Center, Leiden, The Netherlands

⁷ Wellcome Trust Center for Molecular Parasitology, Institute of Infection, Immunity and Inflammation, College of Medical Veterinary & Life Sciences, University of Glasgow, Glasgow, Scotland, United Kingdom

⁸ Antigen Discovery Inc., Irvine, California, USA

^{#a} Current Address: Centre for Molecular Medicine and Therapeutics, Department of Medical Genetics, BC Children's Hospital Research Institute, University of British Columbia, Vancouver, BC, Canada

^{#b} Current Address: Theoretical Biology and Bioinformatics, Department of Biology, Utrecht University, Utrecht, The Netherlands

* Corresponding author

Email Martijn.Huynen@radboudumc.nl (MAH)

[¶] These authors contributed equally to this work.

12 **Abstract**

13

14 *Plasmodium* gametocytes are the sexual forms of the malaria parasite essential for
15 transmission to mosquitoes. To better understand how gametocytes differ from asexual
16 blood-stage parasites, we performed a systematic analysis of available 'omics data for *P.*
17 *falciparum* and other *Plasmodium* species. 18 transcriptomic and proteomic data sets were
18 evaluated for the presence of curated "gold standards" of 41 gametocyte-specific versus 46
19 non-gametocyte genes and integrated using Bayesian probabilities, resulting in gametocyte-
20 specificity scores for all *P. falciparum* genes.

21 To illustrate the utility of the gametocyte score, we explored newly predicted gametocyte-
22 specific genes as potential biomarkers of gametocyte carriage and exposure. We analyzed
23 the humoral immune response in field samples against 30 novel gametocyte-specific
24 antigens and found five antigens to be differentially recognized by gametocyte carriers as
25 compared to malaria-infected individuals without detectable gametocytes. We also validated
26 the gametocyte-specificity of 15 identified gametocyte transcripts on culture material and
27 samples from naturally infected individuals, resulting in eight transcripts that were >1000-fold
28 higher expressed in gametocytes compared to asexual parasites and whose transcript
29 abundance allowed gametocyte detection in naturally infected individuals.

30 Our integrated genome-wide gametocyte-specificity scores provide a comprehensive
31 resource to identify targets and monitor *P. falciparum* gametocytemia.

32

33

34 Introduction

35

36 Despite a decrease in malaria incidence and mortality over the past two decades, malaria
37 remains a major global health challenge ^{1,2}. Furthermore, the emergence and spread of
38 insecticide resistance in mosquitoes ³ and artemisinin resistance in *Plasmodium falciparum*
39 (*Pf*) ⁴⁻⁶ threaten recent gains in malaria control. The decline in malaria burden and the
40 necessity to contain artemisinin-resistance have increased interest in malaria elimination
41 that may require interventions that specifically aim to prevent malaria transmission. Malaria
42 transmission depends on male and female gametocytes, the sexually reproducing forms of
43 the *Plasmodium* parasite that are ingested by blood-feeding *Anopheles* mosquitoes. In the
44 mosquito gut, gametocytes may complete the parasite's reproductive cycle and, following
45 sporogonic development, render the mosquito infectious. Factors that govern gametocyte
46 production and infectivity remain poorly understood. Whilst recent studies have shed light on
47 the processes controlling gametocyte commitment ^{7,8}, commitment and maturation of
48 gametocytes may differ between infections and over the course of infections, under
49 influence of environmental and host factors ^{9,10}. A better understanding of gametocyte
50 dynamics during infections, as well as the development of tools to monitor or target
51 gametocytes, may be informed by high-throughput protein and transcriptome studies ^{11,12}. In
52 the past 15 years, a number of large-scale studies on *Plasmodium* gametocytes have been
53 reported: the proteome of *Pf* and the rodent malaria parasite *Plasmodium berghei* (*Pb*) have
54 been examined by mass spectrometry ¹³⁻²¹, and the transcriptome of both species by
55 microarray and RNA sequencing ^{15,20,22-28}. These studies differed in their focus and
56 resolution in examining (sexual) developmental stages and each faced challenges in
57 detecting low abundance proteins ²⁹ and by the purity of parasite populations ^{13,14}. The use
58 of fluorescent parasites and fluorescence-assisted sorting of staged parasites have recently
59 permitted a better discrimination of proteins in either male or female gametocytes ^{16,17,20} and
60 have allowed more detailed comparisons of *Plasmodium* life-stages. However, individual
61 studies are still vulnerable to imperfect sample purity, and other sources of uncertainty such
62 as correct gene identification for accurate peptide assignment. These technical and
63 methodological challenges lead to discrepancies between individual studies and hamper firm
64 conclusions about gametocyte-specificity of proteins and transcripts.

65

66 We utilized the numerous published proteomics and transcriptomics *Plasmodium* data sets
67 in a comprehensive data integration framework to obtain a consensus of gametocyte-
68 specific transcripts and proteins. Our data integration approach is an adaptation of the naïve

69 Bayesian classifiers that have previously been applied in the prediction of protein
70 interactions and components of cellular systems^{30,31}. The framework calculates probabilities
71 that any given transcript or protein is gametocyte-specific given the evidence presented
72 across the total of transcriptomics and proteomics data. A key aspect of the methodology is
73 that it takes into account the predictive power of each contributing data set: it assigns
74 weights to data sets based on their ability to distinguish gold standard lists of gametocyte
75 and asexual proteins. These we have constructed using existing literature where life-stage
76 specificity was confirmed using classical “non-omics” approaches (e.g. protein detection in
77 immunofluorescence-assays, functional/genetic studies), followed by expert curation. The
78 most informative data (from datasets with the highest discriminative power against the gold-
79 standard lists) will thus contribute most to the predictions, while less informative data are
80 down-weighted. This allows for (i) the resolution of conflicting evidence without disregarding
81 data, and (ii) the construction of a transparent scoring system in which the relative
82 contribution of each data set is directly visible.

83

84 Using this approach, we propose a robust gametocyte-specificity score for all *Pf* genes that
85 allows a consensus list of gametocyte-specific genes at protein and transcript level. We
86 illustrate the utility of our findings by examining naturally acquired responses to newly
87 identified gametocyte-specific proteins in gametocyte-carriers and non-carriers by protein
88 microarray. In addition, we confirmed gametocyte-specificity for a selection of gametocyte-
89 specific transcripts using culture material from geographically distinct *Pf* strains and samples
90 from naturally infected malaria patients.

91

92

93 Results

94

95 Weighted integration of proteomics and gene expression data

96 Using Bayesian statistics, we integrated *Plasmodium* mass spectrometry and transcript
97 datasets from 18 different studies (Table 1) on *P. falciparum* (*Pf*; n=14), *P. berghei* (*Pb*; n=3)
98 and *P. vivax* (*Pv*; n=1). Since gametocyte biology differs between *Plasmodium* species,
99 scores were calculated for the total set of *Plasmodium* studies and for *Pf* only. Unsupervised
100 clustering of genes based on peptide counts or mRNA expression resulted in grouping
101 according to data acquisition method rather than parasite stage (Fig 1), illustrating the
102 necessity of a supervised approach to discriminate between gametocyte-specific and non-
103 gametocyte genes. To objectively assess the value of individual data sets and allow their
104 assembly into a gametocyte-specificity score, we created a gold standard that served as a

105 benchmark for every sample. This gold standard was collected from literature review and
106 comprises two lists; one of asexually expressed proteins, mainly blood stage but also
107 sporozoite and liver stage, and one of known gametocyte proteins (Supplementary Table
108 S1). A gametocyte-specificity score was then derived for each gene by comparing its
109 expression in all studies to the relative expression of the gametocyte and asexual gold
110 standards in those samples (Supplementary Fig S1-2). Proteins or RNAs detected in a study
111 with high discriminative power for gametocyte and asexual gold standard genes (ratio of
112 gametocyte to asexual gold standard genes) received higher gametocyte-specificity scores
113 than those detected in a study with lower discriminative power. The individual log-
114 transformed scores per gene were combined for proteomics and transcriptomics data
115 separately. Scores for *Pf*-only studies (Supplementary Table S2) and all combined data sets
116 were highly correlated (Fig 2D; Pearson's $r=0.9867$ and $r=0.9514$ for proteomics and
117 transcriptomics, respectively) and the latter were used in the remainder of the manuscript.
118 The distribution of scores for proteins and transcripts are presented for the two sets of genes
119 of the gold standard as well as for all other genes (Fig 2A). As expected, the gametocyte
120 and asexual gold standard set of genes are perfectly separated by their respective
121 proteomics-derived scores and show only little overlap in their transcriptomics-derived
122 scores (Fig 2B). The shift in the density peak of the proteomics compared to transcriptomics
123 is due to an inherent property of the method that gives a negative score for gametocyte-
124 specificity to all proteins that were not detected in the proteomics studies ($n=1583$). The
125 highest scoring 100 genes for proteomics and transcriptomics contained 26 (63.4 %) and 15
126 (36.6 %) of the 41 gold standard gametocyte genes, respectively (Fig 2C), indicating both
127 the respective discriminating power of the gold standard and that many other genes are as
128 specific for gametocytes as the highest gold standard representatives.

129 Translationally repressed genes are common in late stage female gametocytes^{32,33} and are
130 detectable by high transcriptomics and low proteomics score. In our analysis 461 genes
131 have this profile (Supplementary Table S3), including genes that are known to be
132 translationally repressed like Pfs28^{32,34} and 186 genes with a previously reported bias
133 towards expression in female gametocytes²⁰.

134

135 **Cross-validation illustrates the improved predictive power of the integrated data.** Ten-
136 fold cross-validation was performed using random subsamples of the gold standard lists to
137 predict the ranks of left-out genes. The resulting proteomics ranking shows near perfect
138 sensitivity, with all but two gold standard gametocyte genes ranking higher than the gold
139 standard asexual genes (Fig 3). The added value of our integrated approach is illustrated by
140 the receiver operating characteristic curve where the integration of data sets gave higher
141 sensitivity and areas under the curve for both proteomics and transcriptomics than any

142 individual study (Fig 3). Using the Bayesian integration based on the complete gold
143 standards, we ranked all *Pf* proteins by giving them a gametocyte-specificity score
144 (Supplementary Table S2). All proteins with a score >5 (n=602) were considered
145 gametocyte-specific. Most of these have not consistently been described as “specific” or
146 “enriched” in gametocytes in the original data sets (Fig 4A and Supplementary Table S4).
147 Previous studies defined 315¹³ to 1725²⁰ proteins as gametocyte-specific for *Pf*. Not only
148 did our integrated approach lead to a better recovery of gold standard listed known
149 gametocyte proteins, we also identified 178 genes with undescribed function as gametocyte-
150 specific (Supplementary Table S2). We further identify a number of proteins as gametocyte-
151 specific even though they had been reported as asexual by previous studies (Fig 4B).
152 A recent proteomics study of male and female *Pf* gametocytes³⁵, not included in our original
153 analysis, was used to test the robustness of our scores. When we included this data set in
154 our final Bayesian proteomics scores, both gametocyte scores and gene ranks before and
155 after addition of this data set were highly correlated (Pearson’s $r=0.997$ and Spearman’s
156 $\rho=0.995$, respectively). Furthermore, the top 100 gametocyte proteins did not change and
157 the top 602 proteins were 96% identical (578 of 602). Taken together, cross-validation and
158 independent data suggest that the integrated gametocyte-specificity score is robust and
159 contains potential novel gametocyte markers.

160

161 **Predicted gametocyte-specific proteins are recognized by gametocyte-carriers.**

162 As an illustration of the utility of gametocyte-specific proteins as markers of gametocyte
163 exposure, we utilized protein microarray data from a study that aimed to characterize the
164 immune profile associated with transmission-reducing immunity in naturally infected
165 gametocyte carriers (Stone, Campo *et al.* 2017 accepted manuscript³⁶). For the current
166 study, we compared responses to our gold standard gametocyte genes (n=40) and novel
167 gametocyte genes from our 100 highest scoring proteins that were on the array (n=30).
168 Antibody prevalences for these genes were compared between Gambian gametocyte
169 carriers and Gambians who carried asexual parasites but not gametocytes as determined by
170 microscopy. Antibody responses to the predicted gametocyte-specific proteins were
171 significantly higher in gametocyte carriers ($p=0.005$), while for the gold standard antigens
172 this difference was less significant ($p=0.058$) (Fig 5A, Mann-Whitney U test). When antigens
173 were analysed individually, a significantly higher antibody prevalence in gametocyte carriers
174 was detected for five novel gametocyte antigens (Fig 5B, $p<0.05$ in Fisher’s exact, corrected
175 for multiple testing, Supplementary Table S5, Supplementary Fig S3). Only two of these five
176 have an assigned function – a DNA ligase, and Gamete egress and sporozoite-traversal
177 protein (GEST). For two of the three remaining *Plasmodium* proteins, we were able to
178 predict a function based on homology, using the sensitive homology detection tool HHpred³⁷.

179 PF3D7_1251000 is homologous to the co-chaperone HSP20 heat shock protein and
180 PF3D7_1439600 is homologous to the MLRQ subunit of complex IV of the oxidative
181 phosphorylation, underlining the enrichment in mitochondrial proteins as discussed below
182 (Supplementary Table S6 includes homology predictions for all conserved, highly
183 gametocyte-specific *Plasmodium* proteins).

184

185 **Gametocyte-specific RNA transcripts detect (submicroscopic) gametocyte carriage**

186 Of the 100 highest-scoring transcripts, 15 non-gold standard candidates were selected for
187 qRT-PCR validation based on their gametocyte scores in a preliminary analysis (Table 2).
188 Mature gametocytes of four *Pf* strains from different geographical origins were compared to
189 asexual blood stage parasites. The minimum transcript abundance difference ($Ct_{\text{Asexuals}} -$
190 $Ct_{\text{Gametocytes}}$) ranged from 4.76 to 14.95 (Fig 6A and Supplementary Table S7 qPCR &
191 primers), reflecting 27.1 to 31,500-fold higher transcript numbers in gametocytes compared
192 to asexual parasites and confirming pronounced upregulation of all selected targets in
193 gametocytes. With a very conservative threshold of 1,000-fold enrichment in three of the four
194 strains tested (Fig 6A), eight of the 15 tested transcripts were highly specific to gametocytes.
195 Transcript abundance in ring-stage parasites was assessed and compared to Pfs25 mRNA,
196 an established and highly abundant yet intron-less female gametocyte specific transcript^{11,38}.
197 Five out of eight gametocyte specific transcripts were undetectable in asexual ring stages at
198 $\leq 10^5$ parasites/mL, similar in specificity to Pfs25 (Fig 6B); the five most sensitive gametocyte
199 markers detected gametocytes across the range of $10^2 - 10^6$ gametocytes/mL (Fig 6C). In
200 RNA samples from a previously reported clinical trial conducted in Kenya³⁹, all eight
201 gametocyte markers detected gametocytes at densities below 10^3 /mL (Fig 6D).

202

203 **Gametocyte-specific proteins are enriched for cytoskeletal movement and** 204 **metabolism functions**

205 To uncover novel characteristics underlying gametocyte function, we analyzed over-
206 represented gene ontology (GO) terms in our integrated consensus gametocyte proteins.
207 The 100 highest ranked proteins were examined for enrichment of GO terms that reflect
208 specific biological functions. Microtubule based movement, metabolism of carboxylic acids
209 and metabolism of nucleic acids were highly enriched among gametocyte proteins
210 (Supplementary Fig S4). Four out of the six putative *Pf* dynein heavy chain proteins are
211 found back among the 100 most gametocyte-specific proteins, alongside a tubulin gamma
212 chain and a tubulin chaperone. The importance of DNA elongation and ligation processes is
213 reflected in GO term associations as well as in antibody response to a DNA ligase (Fig 5B).
214 The “classic” GO term enrichment calculation was complemented by a rank-based gene set

215 enrichment analysis (GSEA). GSEA uses all *Pf* genes and their respective (proteomics
216 derived) gametocyte-specificity score and thus does not include an arbitrary cutoff of the
217 proteins that are or are not gametocyte-specific. It confirmed the above-mentioned results,
218 and in addition to those the terms “mitochondrial protein complex“ (GO:0098798) and “TCA”
219 (GO:0006099) were enriched, stressing both mitochondrial location and processes.
220 Although male gametocytes carry pre-synthesized proteins to rapidly form eight motile
221 gametes upon activation in the mosquito midgut, we did not observe flagellum associated
222 terms in the GSEA. The reason for this is that current GO term annotation for *Pf* has only
223 one gene with the “cilium” GO term for cellular component (GO:0005929; PF3D7_1025500),
224 one for “axoneme” (GO:0005930; PF3D7_0828700) and none with the biological-process
225 terms “cilium or flagellum-dependent cell motility” (GO:0001539) or “axoneme assembly”
226 (GO:0035082). We supplemented the GO annotation with a list of 28 *Pf* cilium genes
227 (Methods). The newly assembled “cilium” GO term now acquired the highest enrichment
228 score in the GSEA (Supplementary Fig S5). This may reflect the formation of the flagella of
229 the microgamete but may also (partially) reflect intracellular trafficking and or be associated
230 with genome replication as the term has overlap with the genes annotated for microtubule
231 processes (Supplementary Table S5 GO terms).
232

233 Discussion

234

235 Combining proteomic and transcriptomic data from 18 sources, we present an integrated
236 consensus score for gametocyte-specific proteins and transcripts. We predict 602
237 gametocyte-enriched proteins of which 186 are currently without ascribed function. We
238 illustrate the potential utility of our gametocyte score by providing evidence for differential
239 recognition of gametocyte proteins by naturally infected gametocyte carriers and the
240 sensitive detection of mRNA of novel gametocyte transcripts in field samples.

241

242 The gametocyte proteome of *P. falciparum* (*Pf*) has been assessed repeatedly. Individual
243 lists of gametocyte-specific proteins^{13,14,16,17,19,20} have unavoidable limitations related to
244 comparator (asexual) parasite stages, sample purity, assay sensitivity and arbitrary cut-offs
245 used to define gametocyte-specificity and show only partial agreement. To acquire a more
246 robust gametocyte-specificity score, we integrated data from these individual studies, along
247 with studies of purified asexual parasites and related *Plasmodium* species. Including gene
248 expression data from multiple species generally increases the likelihood that the combined
249 gene expression data reflect underlying biology, as observed in the *Apicomplexa*^{40,41}. We
250 applied a Bayesian classifier first applied to ‘omics data by Jansen and colleagues³⁰ and

251 adapted by Van der Lee and colleagues to identify genes involved in anti-viral immune
252 responses³¹. The probabilistic approach combines the evidence from all studies in an
253 unbiased way, without giving *a priori* preference of one study over another. Instead, the
254 measurements of all studies were weighted inherently during the scoring process by
255 assessing the retrieval of a gold standard set of genes. As these gametocyte and asexual
256 gold standard sets are of central importance to the study, they have undergone expert
257 curation (see Methods and Acknowledgments). The power of this integrative approach lies
258 not only in weighting data sets by the retrieval of gold standard genes but also in the
259 opportunity to exclude proteins from the gametocyte-specific list by appreciating their
260 presence in (other) asexual samples. A further strength of the approach is that it allows the
261 ranking of gametocyte proteins that have only been reported in a subset of studies. Our
262 integration of data sets reveals that 602 proteins are likely to be specific to gametocytes
263 although very few gametocyte-specific proteins were detected in every underlying dataset
264 and seven proteins had never before been reported as gametocyte-specific.

265 A general limitation of all mass spectrometry (MS) studies is their bias toward highly
266 abundant proteins. Proteins with low-level expression may be missed in a bulk proteome
267 analysis. After integration of the MS studies listed in Table 1, 1583 *Pf* proteins were never
268 detected, representing approximately 28% of all proteins encoded by *Pf*. Some of these
269 might be of too low abundance or expressed during sporozoite or gamete, ookinete and liver
270 stage, which are underrepresented or not included in our data, respectively. New advances
271 in MS that include the sensitive detection of peptides from currently understudied
272 *Plasmodium* life stages may shed light on these currently uncharacterized genes. In addition,
273 approaches that focus specifically on post translational modifications like phosphorylation of
274 proteins as has been done for asexual parasites⁴²⁻⁴⁵ may add new lines of evidence
275 towards gametocyte-specific functions of proteins. Our approach suggests that the currently
276 available MS data is sufficiently comprehensive to identify stage-specific proteins when
277 analysed in an integrative approach. We examined this directly by incorporating a new *Pf*
278 gametocyte MS study³⁵ in our scoring. The authors reported 44 new gametocyte-specific
279 proteins that were not reported by earlier studies. We compared this data set to our
280 integrated data set and found 24 of the 44 had been detected in one or more erythrocytic
281 stages or sporozoites^{14,18,42,46,47} while 11 others had been identified in a (single) gametocyte
282 sample before (Supplementary Fig S6). Importantly, the scores and top 100 gametocyte
283 genes remained unaltered by integrating this new dataset.

284

285 The ranking of gametocyte-specificity that we provide here can i) aid in understanding the
286 biology of this life stage and ii) improve diagnostics related to gametocyte exposure and
287 carriage. Regarding gametocyte biology, our high-ranking gametocyte-specific genes are

288 enriched for mitochondrial, metabolism and microtubule processes and DNA replication,
289 supporting the quality of the data integration. The enrichment of mitochondrial localization
290 and process is consistent with what we know about the enlarged mitochondrion of
291 gametocyte stages⁴⁸ and increased activity of the citric acid cycle⁴⁹. DNA replication terms
292 are highly enriched which is consistent with what happens in the subsequent life stage in
293 which the (micro)gamete rapidly duplicates its genomic DNA three times. Regarding the use
294 of the gametocyte score to inform gametocyte diagnostics, diagnostics can directly detect
295 nucleic acids specific to gametocytes¹¹ or detect antibody responses reflecting past/recent
296 exposure as is increasingly used for asexual *P. falciparum* and *P. vivax* parasites^{50,51}. We
297 use our integrated gametocyte list to explore its utility for both approaches. We validated 15
298 transcript targets in four different *Pf* strains, comparing transcript abundance in gametocytes
299 and asexual parasites. All tested targets were enriched in gametocytes. Five targets were
300 tested for their sensitivity and can recognize 100 gametocytes/mL, while the signal is
301 undetectable when fewer than 10⁵-10⁶ ring-stage parasites/mL are present. In practical
302 terms, these markers may be used to reliably detect gametocytes at densities well below the
303 microscopic threshold of detection in samples without high-densities of asexual parasites,
304 similar to the gametocyte marker that is currently most widely used, the female gametocyte-
305 specific Pfs25³⁸.

306 As an alternative approach to the detection of gametocyte carriage in populations, we
307 utilized a gametocyte-enriched protein microarray (Stone, Campo *et al.* 2017 accepted
308 manuscript³⁶) to determine antibody responses to genes that we here describe as highly
309 gametocyte-specific. The bacterial expression system used for the array has known
310 limitations with the expression of conformational proteins⁵² and should thus be considered a
311 'rule in' rather than 'rule out' approach to immune recognition. Moreover, the array was
312 constructed with the aim of detecting surface proteins or exported proteins whilst our list
313 does not require these characteristics. Only 30 of our top 100 novel gametocyte antigens
314 were thus printed on this array. Antibody responses to five gametocyte proteins were
315 significantly more prevalent in gametocyte carriers than in carriers of the asexual blood
316 stage only. This is the first evidence that antibody responses may be indicative of current
317 gametocyte carriage. Importantly, the dichotomization of gametocyte-exposed and non-
318 exposed individuals was based on a single time-point screening for gametocytes by
319 microscopy. Microscopy has a low sensitivity for detecting gametocytes that commonly
320 circulate at low densities⁵³ and several of the asexual parasite carriers are likely to have had
321 preceding or concurrent low densities of circulating gametocytes. Antibody prevalence in the
322 group classified as gametocyte-negative by microscopy may thus be associated with
323 concurrent low-density gametocytemia and/or long-lived antibody responses acquired
324 following previous gametocyte exposure. The presently analysed samples thus do not allow

325 any conclusions on a possible role of submicroscopic gametocyte densities in boosting or
326 maintaining antibody responses to gametocyte antigens. Refined studies with longitudinal
327 sampling and gametocyte detection by sensitive qRT-PCR methodologies are needed to
328 formally assess antibody kinetics in relation to gametocyte exposure and determine whether
329 recent markers of exposure to blood stage antigens⁵⁰ can be complemented by a set of
330 markers for recent or long-term gametocyte exposure.

331

332 We described the assembly of a curated gold standard set of gametocyte and asexual
333 proteins and used this new resource to rank the likelihood of all *Pf* proteins and transcripts
334 being specific to the gametocyte stage. Data from 18 publicly available studies were
335 integrated to resolve partially conflicting evidence. The resulting consensus lists can be used
336 for guidance of future investigations as we have shown the value of our predictions by in
337 vitro validation.

338

339 **Materials and Methods**

340

341 **Assembly of a gold standard for gametocyte and asexual proteins to weigh whole** 342 **proteome/transcriptome data sets**

343 To build a gold standard against which the performance of individual data sets could be
344 assessed, we identified proteins that are known to be expressed in either asexual parasites
345 (mostly blood stage, also including sporozoites and liver stage) or gametocytes. This list was
346 initially informed by literature review (Supplementary table S1) for expression in the
347 respective stages as detected by immunofluorescence assays and/or western blot,
348 supplemented with *P. falciparum* blood stage or transmission blocking vaccine candidates.
349 This initial list was then communicated with experts (including the authors DAB, PA, FS,
350 CJJ, SMK, TWAK, MM, CD, RS and TB) and edited. If additional proteins were suggested
351 for inclusion in the list, published evidence was requested and examined prior to inclusion of
352 the protein. The final asexual gold standard list contains 46 proteins; the final gametocyte list
353 contains 41 proteins. These gold standard lists (Supplementary Table 1) represent the
354 balance between very strict inclusion criteria and sufficient set size to evaluate the quality of
355 all data sets integrated. We tested for the detection of these proteins or transcripts in the
356 respective samples, using a Bayesian statistics approach that we have successfully applied
357 previously for genes involved in anti-viral immune defense³¹.

358

359 **Data selection and integration**

360 Data sets that measured protein and transcript abundance in *Pf* gametocytes were balanced
361 with data from other life stages and supplemented with studies of the rodent malaria parasite
362 *P. berghei*. One MS study on *P. vivax* was included as it is based on of *ex vivo* blood
363 material as opposed to all other studies that used *in vitro* cultivated parasites. Unique
364 peptide counts were retrieved from plasmoDB (version 28) in which sequenced peptides
365 from the published studies are always mapped to the most recent genome annotation, or
366 supplementary material of the respective studies. Many aspects determine how well a
367 protein is represented in proteomics data that are obtained via MS, like its length or
368 posttranslational modifications. Remapping original MS data to newly annotated genomes
369 improves the quality of the predicted proteins¹⁹. We were however not able to retrieve those
370 data from the studies^{13,14,16} and therefore decided to take those proteins at face value.
371 Notice that also these early studies contribute significantly to our integrated lists.
372 Expression percentiles were retrieved from plasmoDB (version 28) or calculated from raw
373 data in the respective supplementary material. Gametocyte samples were summarized if
374 applicable (using the maximum peptide count/expression percentile of different stages or
375 male and female gametocytes) as were asexual samples, only considering the highest
376 expression in any sample or time point.
377 For MS and transcriptomics data sets, separate scores for gametocyte-specificity of any *Pf*
378 gene have been calculated. In brief, protein or transcript expression has been categorized
379 from absent to high expression levels as given by number of unique peptides or expression
380 percentiles, respectively. For each of the respective bins, a score was calculated depending
381 on the relative retrieval of gametocyte and asexual gold standard genes. The log ratio of
382 these retrieved genes defined the score for all other genes within the same bin. The final
383 gametocyte score calculates as the prior probability of a gene being gametocyte-specific that
384 is updated using the contributions of the data sets:

385

$$GametocyteScore = \log_2 \left(\frac{P_{gct}}{P_{\sim gct}} \right) + \sum_{i=1}^n \log_2 \left(\frac{P(data_i|Gct)}{P(data_i|\sim Gct)} \right) \quad (1)$$

386

$$with \frac{P(data_i|Gct)}{P(data_i|\sim Gct)} = \frac{gametocyteGS_i}{asexualGS_i} \quad (2)$$

387

388 where gametocyteGS and asexualGS are the fractions of retrieved gametocyte and asexual
389 gold standard genes in sample *i*, respectively. We used a pseudocount of 1 if necessary to
390 prevent division by zero if none of the gold standard genes was retrieved in this specific
391 sample and bin. It was assumed that the likelihood of a gene to be either gametocyte or
392 asexual specific is equally high, thus the (log-transformed) prior equals 0 and the final score

393 depends solely on the integrated data. In the selection of a set of proteins that we assigned
394 to be gametocyte-specific we chose a cutoff score of 5.0 (proteomics-derived). The cutoff
395 score of 5.0 can be interpreted as: a gene has to be $2^5 = 32$ times more likely to be
396 gametocyte-specific than asexual specific. The score of 5.0 was based on the behavior of
397 the gold standard genes. Out of the 41 gametocyte gold standard genes, 37 have a score
398 higher than 5.0, while none of the asexual gold standard genes do.
399 When applicable, genes from *Pb* and *Pv* were treated as their respective *Pf* orthologs as
400 retrieved from plasmDB⁵⁴, to be able to integrate all data sets. When no ortholog is known,
401 the respective non-*Pf* data sets did not contribute to the score of this particular gene. Scores
402 using *Pf* data exclusively were also calculated (Supplementary Table S2 includes all scores
403 and rankings with expression information from all integrated studies).

404

405 **Cross-validation of the scoring method**

406 We performed a ten-fold cross-validation to assess the predictive performance of the
407 integrated gametocyte-specificity score (i.e. its ability to discriminate known gametocyte vs.
408 asexual genes). For that, we subsampled both gold standard gene sets ten times (folds),
409 without replacement (i.e. each gene is selected exactly once). Then for each fold we re-
410 weighed and integrated the data sets based on nine-tenth of gold standard genes, and
411 collected the ranks of the one-tenth of genes that were left out in that particular fold. A ROC
412 curve was constructed based on those ranks. Using the same strategy, ROC curves for
413 individual data sets that comprised both gametocyte and asexual samples were constructed
414 for comparison.

415

416 **Protein microarray to measure humoral immune responses**

417 A protein microarray that was enriched for gametocyte proteins was produced and probed
418 as described earlier for a study aiming to unravel the immune signature of naturally acquired
419 transmission-reducing immune responses in gametocyte carriers (Stone, Campo *et al.* 2017
420 accepted manuscript³⁶). As a control group, Gambian asexual parasite carriers without
421 gametocytes detectable by microscopy were included in the probing. For the current study
422 array data from this control group (n=63) and Gambian gametocyte carriers were used
423 (n=164)⁵⁵⁻⁶⁰. All of these 227 individuals were sampled during a period of intense malaria
424 transmission intensity in The Gambia and likely had (multiple) previous malaria infections
425^{55,61}. For these populations, responses to 30 newly defined highly gametocyte-specific
426 antigens (from 24 genes) were compared between gametocyte carriers and non-carriers
427 (Mann-Whitney U test). Seropositivity for each of the antigens was determined using a
428 mixture model-based cutoff and related to gametocyte carriage using Fisher's exact test,

429 corrected for multiple testing (Benjamini-Hochberg) of a total of 70 antigens (including 40
430 antigens from the gametocyte gold standard).

431

432 **Transcript abundance in different life stages and strains**

433 The abundance of 15 predicted gametocyte-specific targets was measured in asexual
434 parasites and gametocytes of four different *Pf* strains from *in vitro* culture. The targets were
435 selected from the 100 highest scoring transcripts to account for uncertainties about the
436 absolute scoring of transcriptomics data with a protein-based gold standard. We do not
437 assume a clear hierarchy between these top 100 scoring transcripts and consider any of
438 these genes highly gametocyte-specific. The 15 highest-ranking non-gold standard genes
439 were selected based on a preliminary analysis of the data, and contain genes that are
440 currently not annotated as well as genes with known protein function in gametocytes (PUF1
441⁶² and Ccp4⁶³). In the final generation of the gametocyte-scores, all validation genes were
442 retained in the top 100 scoring genes. The *Pf* strains used are of West African (NF54,
443 NF166, NF175) and Southeast Asian origin (NF135). All strains were cultured and
444 synchronized as described previously²⁰. Using established standard curves, the same
445 concentrations of parasites were compared for Ct values in qRT-PCR (for primers, see
446 Supplementary Table S7). Extracted nucleic acids were DNase-treated before reverse
447 transcription when introns were absent from the targets. Initial comparison was between
448 mixed asexual blood stage parasites (considering the lowest Ct measured in any strain and
449 replicate) and stage V gametocyte (highest Ct measured per strain and replicate). Promising
450 targets with a high $Ct_{\text{Asexuals}} - Ct_{\text{Gametocytes}}$ were further examined in serial dilutions of stage V
451 gametocytes and synchronized asexual material of the strain NF54 (10, 20, 30, 40 hours
452 post invasion, resembling early rings, late rings, trophozoites and schizonts, respectively).
453 All qRT-PCR reactions were analyzed in technical triplicates, from biological triplicates
454 (NF54) or duplicates (remaining strains).

455 RNA samples were used from a clinical malaria trial conducted in Western Kenya³⁹.

456 Samples from days 3 and 7 after treatment were selected to ensure a range of (low-density)
457 gametocyte carriage to test qRT-PCR sensitivity.

458

459 **Go term enrichment in top 100 proteins or rank-based enrichment**

460 Current GO term annotation for *Pf* was retrieved from plasmoDB (release 30) and analyzed
461 using the topGO R package⁶⁴ for the enrichment of terms in the 100 highest scoring
462 proteins versus all *Pf* proteins. Semantic clustering of the significant GO terms in Biological
463 Process ontology was done with the Revigo webtool for *Pf*⁶⁵. Second, gene set enrichment
464 analysis (GSEA) based on all *Pf* proteins and their ranks and scores was performed using
465 the software available at <http://software.broadinstitute.org/gsea/downloads.jsp>. Based on the

466 cilium genes reported by the Syscilia consortium⁶⁶, we assembled a “cilium” GO term of
467 mixed ontology (GO:9999999) for *Pf* with 28 predicted orthologs (Supplementary Table S6).

468

469 **Data availability**

470 All data generated or analysed during this study are included in this published article (and its
471 Supplementary Information files).

- 472 - The gold standard lists (Supplementary Table S1)
- 473 - Bayesian gametocyte scoring for proteomics and transcriptomics data. Includes *Pf*
474 only-scores and expression values for any gene and individual data set
475 (Supplementary Table S2)
- 476 - Potential translationally repressed genes with high transcriptomics score (>7) and
477 low proteomics score (<-10, Supplementary Table S3)
- 478 - Overview over previously reported gametocyte-specificity per study (Supplementary
479 Table S4)
- 480 - Seroprevalence in The Gambia in gametocyte carriers and non-carriers
481 (Supplementary Table S5)
- 482 - Function predictions for highly gametocyte-specific proteins with lacking annotation
483 (Supplementary Table S6)
- 484 - Transcript validation in 15 targets, including primer sequences for qRT-PCR
485 (Supplementary Table S7)
- 486 - GO term analyses with cilium genes (Supplementary Table S8)

487

488

489 References

490

- 491 1. WHO. Reversing the Incidence of Malaria 2000–2015. in (2015).
- 492 2. Bhatt, S. *et al.* The effect of malaria control on *Plasmodium falciparum* in Africa between
493 2000 and 2015. *Nature* **526**, 207–11 (2015).
- 494 3. Ranson, H. & Lissenden, N. Insecticide Resistance in African Anopheles Mosquitoes: A
495 Worsening Situation that Needs Urgent Action to Maintain Malaria Control. *Trends Parasitol.*
496 **32**, 187–196 (2016).
- 497 4. Ashley, E. A. *et al.* Spread of Artemisinin Resistance in *Plasmodium falciparum* Malaria. *N.*
498 *Engl. J. Med.* **371**, 411–423 (2014).
- 499 5. Takala-Harrison, S. *et al.* Independent emergence of artemisinin resistance mutations among
500 *Plasmodium falciparum* in Southeast Asia. *J. Infect. Dis.* **211**, 670–679 (2015).
- 501 6. Ménard, D. *et al.* A Worldwide Map of *Plasmodium falciparum* K13-Propeller Polymorphisms.
502 *N. Engl. J. Med.* **374**, 2453–2464 (2016).
- 503 7. Kafsack, B. F. C. *et al.* A transcriptional switch underlies commitment to sexual development
504 in malaria parasites. *Nature* **507**, 248–52 (2014).
- 505 8. Pelle, K. G. *et al.* Transcriptional profiling defines dynamics of parasite tissue sequestration
506 during malaria infection. *Genome Med.* **7**, 19 (2015).
- 507 9. Johnston, G. L., Smith, D. L. & Fidock, D. A. Malaria’s Missing Number: Calculating the Human
508 Component of R0 by a Within-Host Mechanistic Model of *Plasmodium falciparum* Infection
509 and Transmission. *PLoS Comput. Biol.* (2013). doi:10.1371/journal.pcbi.1003025
- 510 10. Bousema, T. & Drakeley, C. Epidemiology and infectivity of *Plasmodium falciparum* and
511 *Plasmodium vivax* gametocytes in relation to malaria control and elimination. *Clinical*
512 *Microbiology Reviews* **24**, 377–410 (2011).
- 513 11. Joice, R. *et al.* Inferring Developmental Stage Composition from Gene Expression in Human
514 Malaria. *PLoS Comput. Biol.* **9**, 1–13 (2013).
- 515 12. Proietti, C. & Doolan, D. L. The case for a rational genome-based vaccine against malaria.
516 *Frontiers in Microbiology* **6**, (2015).
- 517 13. Lasonder, E. *et al.* Analysis of the *Plasmodium falciparum* proteome by high-accuracy mass
518 spectrometry. *Nature* **419**, 537–42 (2002).
- 519 14. Florens, L. *et al.* A proteomic view of the *Plasmodium falciparum* life cycle. *Nature* **419**, 520–6
520 (2002).
- 521 15. Hall, N., Karras, M. & Raine, J. D. A Comprehensive Survey of the *Plasmodium* Life Cycle by
522 Genomic, Transcriptomic, and Proteomic Analyses. *Science (80-.)*. **307**, 82–86 (2005).
- 523 16. Khan, S. M. *et al.* Proteome Analysis of Separated Male and Female Gametocytes Reveals
524 Novel Sex-Specific *Plasmodium* Biology. *Cell* **121**, 675–687 (2005).
- 525 17. Silvestrini, F. *et al.* Protein Export Marks the Early Phase of Gametocytogenesis of the Human
526 Malaria Parasite *Plasmodium falciparum*. *Mol. Cell. Proteomics* **9**, 1437–1448 (2010).
- 527 18. Oehring, S. C. *et al.* Organellar proteomics reveals hundreds of novel nuclear proteins in the
528 malaria parasite *Plasmodium falciparum*. *Genome Biol* **13**, R108 (2012).
- 529 19. Tao, D. *et al.* Sex-partitioning of the *Plasmodium falciparum* Stage V Gametocyte Proteome
530 Provides Insight into *falciparum*-specific Cell Biology. *Mol. Cell. Proteomics* **13**, 2705–2724
531 (2014).
- 532 20. Lasonder, E. *et al.* Integrated transcriptomic and proteomic analyses of *P. falciparum*
533 gametocytes: molecular insight into sex-specific processes and translational repression.
534 *Nucleic Acids Res.* **44**, gkw536 (2016).
- 535 21. Suárez-Cortés, P. *et al.* Comparative proteomics and functional analysis reveal a role of *P.*
536 *falciparum* osmiophilic bodies in malaria parasite transmission. *Mol. Cell. Proteomics* (2016).
537 doi:10.1074/mcp.M116.060681
- 538 22. Bozdech, Z. *et al.* The transcriptome of the intraerythrocytic developmental cycle of
539 *Plasmodium falciparum*. *PLoS Biol.* **1**, 85–100 (2003).

- 540 23. Le Roch, K. G. *et al.* Discovery of gene function by expression profiling of the malaria parasite
541 life cycle. *Science (80-.)*. **301**, 1503–1508 (2003).
- 542 24. Young, J. A. *et al.* The Plasmodium falciparum sexual development transcriptome: A
543 microarray analysis using ontology-based pattern identification. *Mol. Biochem. Parasitol.* **143**,
544 67–79 (2005).
- 545 25. Llinas, M. Comparative whole genome transcriptome analysis of three Plasmodium
546 falciparum strains. *Nucleic Acids Res.* **34**, 1166–1173 (2006).
- 547 26. Otto, T. D. *et al.* New insights into the blood-stage transcriptome of Plasmodium falciparum
548 using RNA-Seq. *Mol. Microbiol.* **76**, 12–24 (2010).
- 549 27. López-Barragán, M. J. *et al.* Directional gene expression and antisense transcripts in sexual
550 and asexual stages of Plasmodium falciparum. *BMC Genomics* **12**, 587 (2011).
- 551 28. Otto, T. D. *et al.* A comprehensive evaluation of rodent malaria parasite genomes and gene
552 expression. *BMC Biol.* **12**, 86 (2014).
- 553 29. Wasinger, V. C., Zeng, M. & Yau, Y. Current status and advances in quantitative proteomic
554 mass spectrometry. *Int. J. Proteomics* **2013**, 180605 (2013).
- 555 30. Jansen, R. *et al.* A Bayesian networks approach for predicting protein-protein interactions
556 from genomic data. *Science (80-.)*. (2003). doi:10.1126/science.1087361
- 557 31. van der Lee, R. *et al.* Integrative Genomics-Based Discovery of Novel Regulators of the Innate
558 Antiviral Response. *PLoS Comput. Biol.* **11**, (2015).
- 559 32. Mair, G. R. *et al.* Regulation of Sexual Development of Plasmodium by Translational
560 Repression. *Science (80-.)*. **313**, 667–669 (2006).
- 561 33. Mair, G. R. *et al.* Universal features of post-transcriptional gene regulation are critical for
562 Plasmodium zygote development. *PLoS Pathog.* **6**, (2010).
- 563 34. Miao, J. *et al.* Puf Mediates Translation Repression of Transmission-Blocking Vaccine
564 Candidates in Malaria Parasites. *PLoS Pathog.* (2013). doi:10.1371/journal.ppat.1003268
- 565 35. Miao, J. *et al.* Sex-Specific Biology of the Human Malaria Parasite Revealed from the
566 Proteomes of Mature Male and Female Gametocytes. *Mol. Cell. Proteomics* (2017).
- 567 36. Stone, W. J. R. *et al.* Unravelling the immune signature of Plasmodium falciparum
568 transmission reducing immunity. *Nat. Commun.* (2017).
- 569 37. Söding, J., Biegert, A. & Lupas, A. N. The HHpred interactive server for protein homology
570 detection and structure prediction. *Nucleic Acids Res.* **33**, 244–248 (2005).
- 571 38. Stone, W. *et al.* A Molecular Assay to Quantify Male and Female Plasmodium falciparum
572 Gametocytes: Results From 2 Randomized Controlled Trials Using Primaquine for Gametocyte
573 Clearance. *J. Infect. Dis.* **216**, 457–467 (2017).
- 574 39. Graves, P. M., Gelband, H. & Garner, P. Primaquine or other 8-aminoquinoline for reducing
575 Plasmodium falciparum transmission. *Cochrane database Syst. Rev.* (2015).
576 doi:10.1002/14651858.CD008152.pub4
- 577 40. Huynen, M. A., Snel, B. & Van Noort, V. Comparative genomics for reliable protein-function
578 prediction from genomic data. *Trends in Genetics* (2004). doi:10.1016/j.tig.2004.06.003
- 579 41. Butler, C. L. *et al.* Identifying novel cell cycle proteins in apicomplexa parasites through co-
580 expression decision analysis. *PLoS One* (2014). doi:10.1371/journal.pone.0097625
- 581 42. Treeck, M., Sanders, J. L. L., Elias, J. E. E. & Boothroyd, J. C. C. The Phosphoproteomes of
582 Plasmodium falciparum and Toxoplasma gondii Reveal Unusual Adaptations Within and
583 Beyond the Parasites' Boundaries. *Cell Host Microbe* **10**, 410–419 (2011).
- 584 43. Lasonder, E. *et al.* The plasmodium falciparum schizont phosphoproteome reveals extensive
585 phosphatidylinositol and cAMP-protein kinase A signaling. *J. Proteome Res.* **11**, 5323–5337
586 (2012).
- 587 44. Pease, B. N. *et al.* Global analysis of protein expression and phosphorylation of three stages
588 of plasmodium falciparum intraerythrocytic development. *J. Proteome Res.* **12**, 4028–4045
589 (2013).
- 590 45. Lasonder, E., Green, J. L., Grainger, M., Langsley, G. & Holder, A. a. Extensive differential

- 591 protein phosphorylation as intraerythrocytic *Plasmodium falciparum* schizonts develop into
592 extracellular invasive merozoites. *Proteomics n/a-n/a* (2015). doi:10.1002/pmic.201400508
- 593 46. Bowyer, P. W., Simon, G. M., Cravatt, B. F. & Bogyo, M. Global profiling of proteolysis during
594 rupture of *Plasmodium falciparum* from the host erythrocyte. *Mol Cell Proteomics* **10**, M110
595 001636 (2011).
- 596 47. Lindner, S. E. *et al.* Total and putative surface proteomics of malaria parasite salivary gland
597 sporozoites. *Mol. Cell. Proteomics* **12**, 1127–43 (2013).
- 598 48. Okamoto, N., Spurck, T. P., Goodman, C. D. & McFadden, G. I. Apicoplast and mitochondrion
599 in gametocytogenesis of *Plasmodium falciparum*. *Eukaryot. Cell* **8**, 128–132 (2009).
- 600 49. MacRae, J. I. *et al.* Mitochondrial metabolism of sexual and asexual blood stages of the
601 malaria parasite *Plasmodium falciparum*. *BMC Biol.* **11**, 67 (2013).
- 602 50. Helb, D. A. *et al.* Novel serologic biomarkers provide accurate estimates of recent
603 *Plasmodium falciparum* exposure for individuals and communities. *Proc. Natl. Acad. Sci. U. S.*
604 *A.* 1501705112- (2015). doi:10.1073/pnas.1501705112
- 605 51. Cutts, J. C. *et al.* Immunological markers of *Plasmodium vivax* exposure and immunity: a
606 systematic review and meta-analysis. *BMC Med.* (2014). doi:10.1186/s12916-014-0150-1
- 607 52. Crompton, P. D. *et al.* A prospective analysis of the Ab response to *Plasmodium falciparum*
608 before and after a malaria season by protein microarray. *Proc. Natl. Acad. Sci. U. S. A.* **107**,
609 6958–6963 (2010).
- 610 53. Schneider, P. *et al.* Submicroscopic *Plasmodium falciparum* gametocyte densities frequently
611 result in mosquito infection. *Am. J. Trop. Med. Hyg.* **76**, 470–474 (2007).
- 612 54. Aurrecochea, C. *et al.* PlasmoDB: A functional genomic database for malaria parasites.
613 *Nucleic Acids Res.* (2009). doi:10.1093/nar/gkn814
- 614 55. Drakeley, C. J., Secka, I., Correa, S., Greenwood, B. M. & Targett, G. A. T. Host haematological
615 factors influencing the transmission of *Plasmodium falciparum* gametocytes to *Anopheles*
616 *gambiae* s.s. mosquitoes. *Trop. Med. Int. Heal.* (1999). doi:10.1046/j.1365-3156.1999.00361.x
- 617 56. Targett, G. *et al.* Artesunate reduces but does not prevent posttreatment transmission of
618 *Plasmodium falciparum* to *Anopheles gambiae*. *J. Infect. Dis.* **183**, 1254–9 (2001).
- 619 57. Drakeley, C. J. *et al.* Parasite infectivity and immunity to *Plasmodium falciparum* gametocytes
620 in Gambian children. *Parasite Immunol.* (2004). doi:10.1111/j.0141-9838.2004.00696.x
- 621 58. Sutherland, C. J. *et al.* Reduction of malaria transmission to *Anopheles* mosquitoes with a six-
622 dose regimen of co-artemether. *PLoS Med.* **2**, 0338–0346 (2005).
- 623 59. Dunyo, S. *et al.* Gametocytaemia after drug treatment of asymptomatic *Plasmodium*
624 *falciparum*. *PLoS.Clin. Trials* (2006).
- 625 60. Hallett, R. L. *et al.* Chloroquine/Sulphadoxine-Pyrimethamine for Gambian Children with
626 Malaria: Transmission to Mosquitoes of Multidrug-Resistant *Plasmodium falciparum*. *PLoS*
627 *Clin. Trials* (2006). doi:10.1371/journal.pctr.0010015
- 628 61. Ceesay, S. J. *et al.* Changes in malaria indices between 1999 and 2007 in The Gambia: a
629 retrospective analysis. *Lancet* (2008). doi:10.1016/S0140-6736(08)61654-2
- 630 62. Shrestha, S., Li, X., Ning, G., Miao, J. & Cui, L. The RNA-binding protein PfPuf1 functions in the
631 maintenance of gametocytes in *Plasmodium falciparum*. *J. Cell Sci.* jcs.186908 (2016).
632 doi:10.1242/jcs.186908
- 633 63. Simon, N. *et al.* Sexual Stage Adhesion Proteins Form Multi-protein Complexes in the Malaria
634 Parasite *Plasmodium falciparum*. *J. Biol. Chem.* (2009). doi:10.1074/jbc.M808472200
- 635 64. Alexa A and Rahnenfuhrer J. topGO: Enrichment Analysis for Gene Ontology. *R package*
636 *version 2.24.0.* (2016).
- 637 65. Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. Revigo summarizes and visualizes long lists of
638 gene ontology terms. *PLoS One* **6**, (2011).
- 639 66. van Dam, T. J., Wheway, G., Slaats, G. G., Huynen, M. A. & Giles, R. H. The SYSCILIA gold
640 standard (SCGSv1) of known ciliary components and its applications within a systems biology
641 consortium. *Cilia* **2**, 7 (2013).

- 642 67. Moreno-Pérez, D. A., Dégano, R., Ibarrola, N., Muro, A. & Patarroyo, M. A. Determining the
643 Plasmodium vivax VCG-1 strain blood stage proteome. *J. Proteomics* **113**, 268–280 (2015).
644
645

646 **Acknowledgments**

647

648 We are thankful to Colin Sutherland (LSHTM London, UK) for commenting on the
649 manuscript and providing serum samples that we used to study antibody responses in
650 Gambian individuals. David Conway and Johannes Dessens (LSHTM London, UK) gave
651 helpful feedback on the initial gold standard lists. We further thank Adam D. Shandling,
652 Jozelyn V. Pablo and Andy A. Teng from Antigen Discovery Inc. (Irvine, CA) for their work
653 on the protein microarray.
654 TB was supported by the Netherlands Organization for Scientific Research (nwo.nl), through
655 a VIDI fellowship 016.158.306. The Radboud Institute for Health Sciences (rihs.nl),
656 supported LM through grant R-2765. The Virgo consortium (virgo.nl), grant FES0908,
657 supported the work of RvdL and TJPvD. TWAK is supported by the Netherlands
658 Organization for Scientific Research (NWO-VIDI 864.13.009). The funders had no role in
659 study design, data collection and analysis, decision to publish, or preparation of the
660 manuscript.
661

662 **Author Contributions Statement**

663

664 MAH, TB, RvdL and LMK conceptualised the work. LMK, RvdL, TJPvD, MAH and TB
665 analysed the data. LMK, TB, DAB, PA, FS, CJJ, SMK, TWAK, MM and RS assembled the
666 Gold Standard. KL, MvdVB, WG, RSS, LMK and WS conducted and analysed experiments
667 with samples and resources provided by CD and JJC. LMK wrote the first draft of the
668 manuscript, all authors reviewed the manuscript.

669

670

671 **Competing financial interests**

672 JJC is employed by Antigen Discovery, Inc. The authors declare no further competing
673 interests.

674

675

676

677

678 **Figure Legends**

679

680

681 *Figure 1* **Figure 1. Clustered data sets used in this study with genes ranked according to**
682 **their protein or transcript expression.** Level of expression as detected in the respective
683 samples with unique peptide counts for MS data and percentiles for transcriptomics. The
684 studies are clustered using complete linkage according to their overall gene expression
685 similarities (Euclidean distance). See Table 1 for study keys. Distribution of
686 asexual(a)/gametocyte(g) samples (red/blue) is shown in top bar, proteomics (P) and
687 transcriptomics (T) (dark/light grey) in lower bar.

688

689 *Figure 2* **Figure 2. Gametocyte-specificity scores for *P. falciparum* genes derived from**
690 **proteomics (P) and transcriptomics (T) data sets.** (A) Boxplot for integrated scores for
691 the two gold standard sets and all other Pf genes, derived from proteomics, transcriptomics
692 or all data sets (combined). (B) Density of P and T gametocyte scores, individual gold
693 standard genes and their scores are indicated at the bottom (red, asexual, blue gametocyte).
694 (C) 100 highest ranking proteins and transcripts, gametocyte gold standard in blue. (D)
695 Correlation of the gametocyte-specificity scores derived from all integrated MS studies and
696 Pf MS studies only.

697

698 *Figure 3* **Figure 3. Validation of Bayesian gametocyte scoring with area under the curve**
699 **(AUC) values.** Integrated data and individual data sets are compared by 10-fold cross-
700 validation (subsampling of gametocyte and asexual gold standard sets). Integrated
701 proteomics (P) and transcriptomics (T) scores in bold lines. *P. berghei* data sets in shades of
702 red, individual proteomics and transcriptomics studies with short and long dashes,
703 respectively. See Table 1 for study keys.

704

705 *Figure 4* **Figure 4. Comparison of reported gametocyte-specific proteins in mass**
706 **spectrometry studies.** (A) Proteins reported as gametocyte-specific by six individual
707 studies, agreements on gametocyte-specificity are summarized in the table. Bayesian:
708 gametocyte-specific proteins (n=602) that have a score > 5 after data integration. The
709 overlap with previously published data sets is shown, but not to scale. Overlap between the
710 individual studies is not shown for better visibility. Note that the Lasonder 2002 study
711 includes proteins that were found in gametocytes or gametocytes and gametes. (B) Proteins
712 that were reported as non-gametocytic and are (partially) included after data integration

713

714 *Figure 5* **Figure 5. Seroprevalence in two cohorts of parasite carriers in The Gambia.**
715 (A+B) Antibodies against the highest scoring gametocyte-specific proteins were measured
716 on protein microarrays. Comparison of positivity (mixture-model cutoff) in gametocyte
717 carriers (n=164) and non-carriers (n=63). Gametocyte presence determined by microscopy.
718 All individuals were positive for asexual parasites. (A) Prevalence of antigens from the gold
719 standard (n=40) and predicted gametocyte-specific proteins (n=30), Mann-Whitney U test
720 (B) Antigens of five predicted gametocyte-specific proteins are preferentially recognized by

721 gametocyte carriers. Error bars indicate the upper limit of the 95% confidence interval
722 around the proportion. $p < 0.05$ Fisher's exact test, corrected for multiple testing of a total of
723 70 antigens (Benjamini-Hochberg)

724

725 **Figure 6. Validation of gametocyte-specific targets in qRT-PCR** Targets are sorted
726 for decreasing gametocyte-specificity in all panels, see Table 2. (A) Minimum transcript
727 abundance in blood stage versus gametocytes in different Pf strains. 1000-fold enrichment
728 of transcript in gametocytes over asexuals was assumed when delta-Ct was 10 or higher
729 (dashed line), considering the lowest Ct value detected in any asexual concentration-
730 matched sample. This threshold was not met by the transcripts with gene IDs in grey. (B)
731 Detection limit of eight validated targets alongside Pfs25 in serial dilutions of Pf NF54
732 asexual stage parasites (ring stage parasites 10-20 hours post invasion). (C) Detection limit
733 of the most sensitive targets in serial dilutions of stage V gametocytes. (A-C) For Pf NF54,
734 all $n=3$, other strains $n=2$ biological replicates (error bars: standard error of the mean), all
735 measurements in triplicates. (D) Sensitivity of eight validated targets in Kenyan blood
736 samples of varying gametocyte densities.

737

738 **Table 1. Data sets used for integration**

739 Life stages Spz sporozoites, Ri rings, Troph trophozoites, Schiz schizonts, Mer
 740 merozoites, Gct gametocytes. MS mass spectrometry (highest unique peptide count
 741 in any of the samples), T transcriptomics (highest percentile in any of the samples).
 742 Data from Miao et al. 2017 was integrated after analyses of high scoring proteins,
 743 ranks and scores are included in Supplementary Table 2.
 744 *Asexual microarray data by Llinas and others retrieved from plasmoDB version 28
 745 (Data set “Pfal3D7 real-time transcription and decay”), no accompanying publication.

Study	Reference	Species	Life Stage	Integrated Data
FI02	Florens et al. 2002 ¹⁴	<i>Pf</i>	Spz, Troph, Mer, Gct	MS asexual/Gct
La02	Lasonder et al. 2002 ¹³	<i>Pf</i>	Troph, Schiz, Gct	MS asexual/Gct
Le03	LeRoch et al. 2003 ²³	<i>Pf</i>	Spz, Ri, Troph, Schiz, Mer, Gct	T asexual/Gct
Ha05	Hall et al. 2005 ¹⁵	<i>Pb</i>	Ri, Troph, Schiz, Gct	MS + T asexual/Gct
Kh05	Khan et al. 2005 ¹⁶	<i>Pb</i>	Mixed blood stage, Gct	MS asexual/Gct
Yo05	Young et al. 2005 ²⁴	<i>Pf</i>	Gct	T Gct
LI06	Llinas et al. 2006 ²⁵	<i>Pf</i>	All blood stages, synchronized	T asexual
Ot10	Otto et al. 2010 ²⁶	<i>Pf</i>	All blood stages, synchronized	T asexual
Bo11	Treeck et al. 2011 ⁴²	<i>Pf</i>	Schiz	MS asexual
Lo11	Lopez-Barragan et al. ²⁷	<i>Pf</i>	Troph, Schiz, Gct	T asexual/Gct
Oe12	Oehring et al. 2012 ¹⁸	<i>Pf</i>	Ri, Troph, Schiz	MS asexual
Mo14	Moreno-Perez et al. 2014 ⁶⁷	<i>Pv</i>	Ri, Troph, Schiz	MS asexual
Ot14	Otto et al. 2014 ²⁸	<i>Pb</i>	Ri, Troph, Schiz, Gct	T asexual/Gct
Si14	Silvestrini et al. 2010, Tao et al. 2014 (re-analyzed) ^{17,19}	<i>Pf</i>	Troph, Schiz, Gct	MS asexual/Gct
Ta14	Tao et al. 2014 ¹⁹	<i>Pf</i>	Gct	MS Gct
LI15	Llinas et al. 2015*	<i>Pf</i>	All blood stages, synchronized	T asexual
La16	Lasonder et al. 2016 ²⁰	<i>Pf</i>	Gct	MS + T Gct
Su16	Suarez-Cortes et al. 2016 ²¹	<i>Pf</i>	Gct	MS Gct
Mi17	Miao et al. 2017 ³⁵	<i>Pf</i>	Gct	MS Gct

746

747

748

Rank	Gene ID	Description	Name	Intron-spanning	Least Ct difference
3	PF3D7_1143600	conserved Plasmodium protein, unknown function	—	no	9.37
5	PF3D7_1147200	tubulin--tyrosine ligase, putative	—	no	8.52
6	PF3D7_1026100	conserved Plasmodium protein, unknown function	—	yes	12.94
7	PF3D7_1438800	conserved Plasmodium protein, unknown function	—	yes	8.57
8	PF3D7_0625100	sphingomyelin synthase 2, putative	SMS2	no	11.04
12	PF3D7_0930000	procollagen lysine 5-dioxygenase, putative	—	no	11.96
14	PF3D7_0518700	mRNA-binding protein PUF1	PUF1	yes	8.31
15	PF3D7_0303900	phosphatidylethanolamine-binding protein, putative	—	yes	10.87
16	PF3D7_1466600	conserved Plasmodium protein, unknown function	—	no	5.67
17	PF3D7_1107900	mechanosensitive ion channel protein, putative	MSCS	no	6.33
18	PF3D7_1214500	conserved Plasmodium protein, unknown function	—	yes	10.61
24	PF3D7_1131500	conserved Plasmodium protein, unknown function	—	no	10.41
51	PF3D7_0929600	G2 protein, putative	—	yes	8.37
61	PF3D7_0816800	meiotic recombination protein DMC1, putative	DMC1	yes	13.29
75	PF3D7_0903800	LCCL domain-containing protein	CCp4	yes	12.70

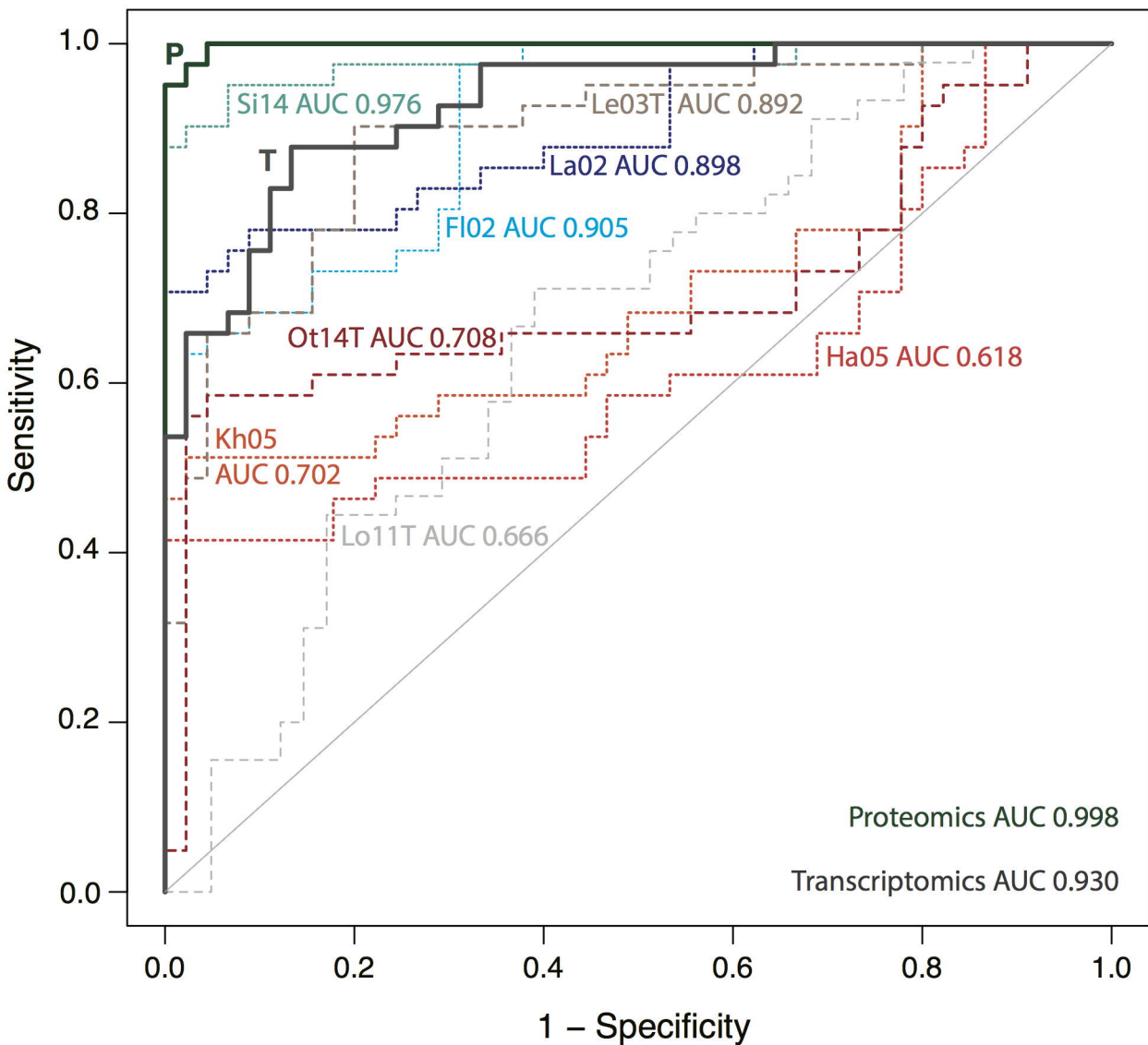
749

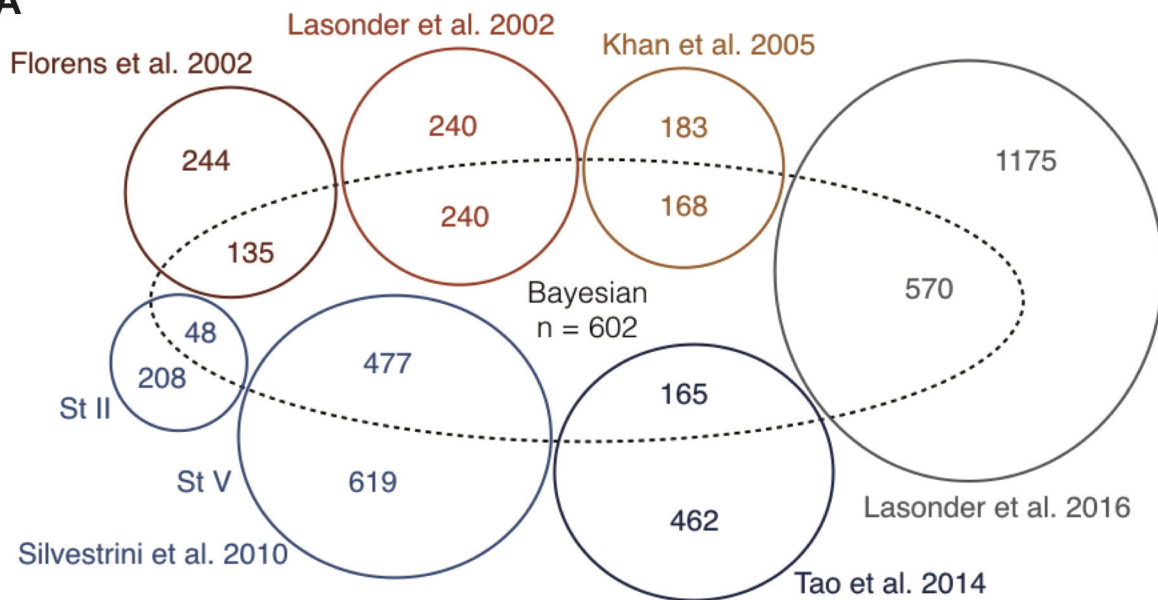
750 **Table 2: Properties of putative gametocyte-specific targets.**

751 Rank in transcriptomics (all data sets) for specificity in gametocytes. Random sample of top
 752 100, excluding the gold standard. If primers are not intron-spanning, samples were DNase I
 753 treated. Ct difference is the difference between the lowest Ct detected in asexual samples
 754 and the highest Ct in concentration-matched stage V gametocytes, averaged across strains
 755 *Pf* NF54, NF135, NF166 and NF175.

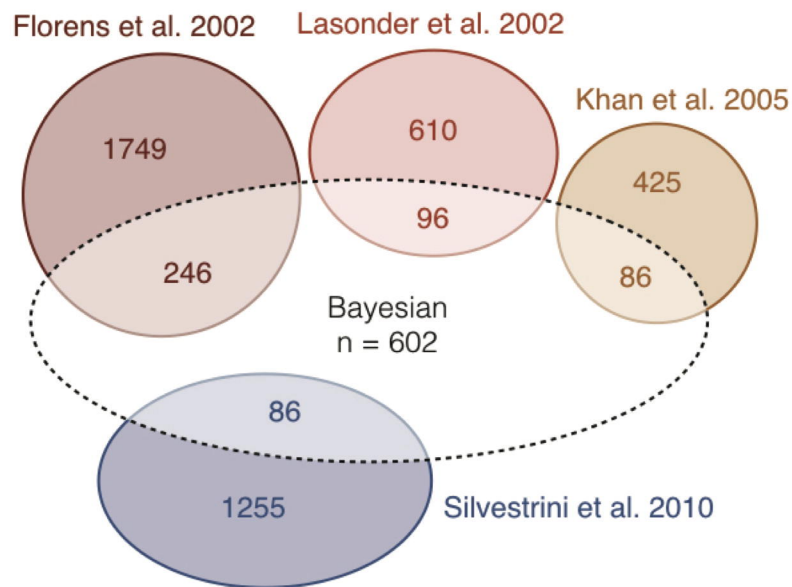
756

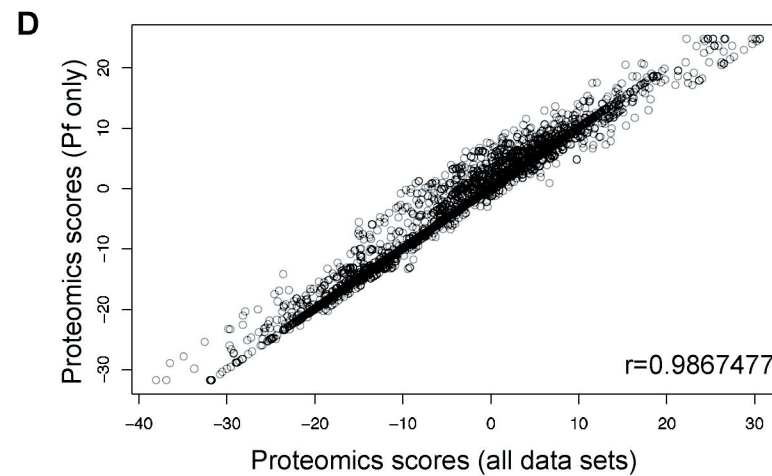
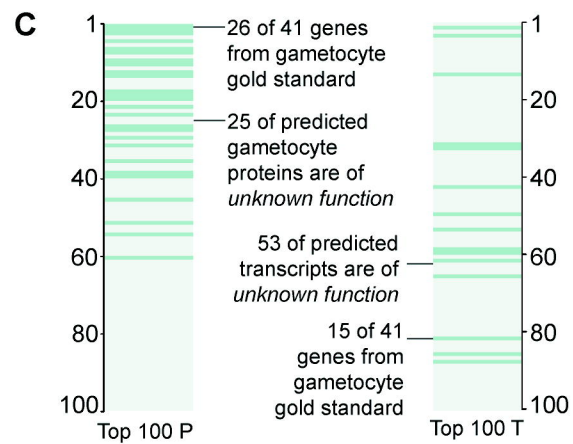
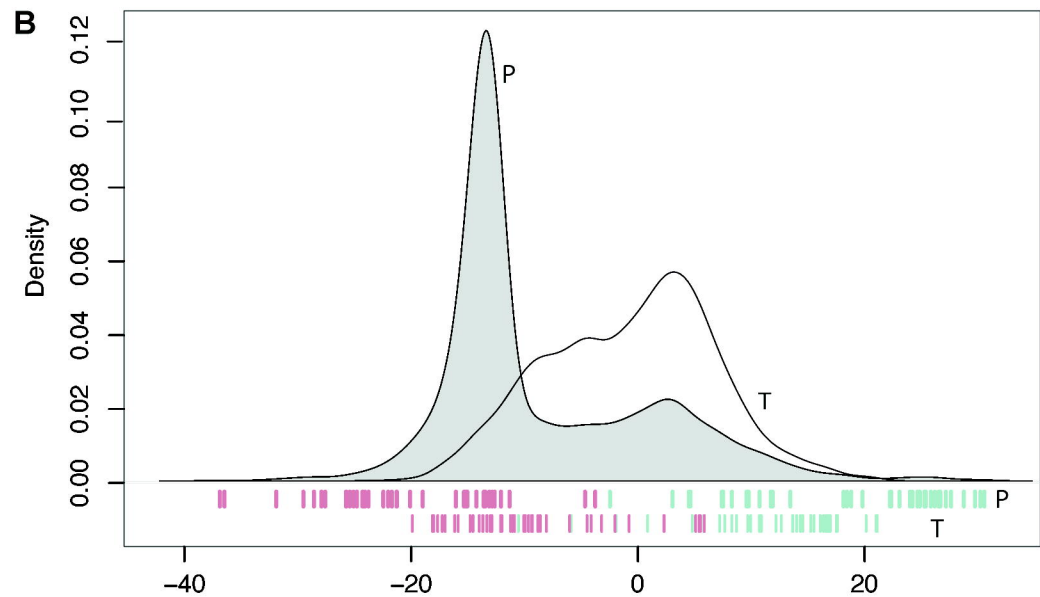
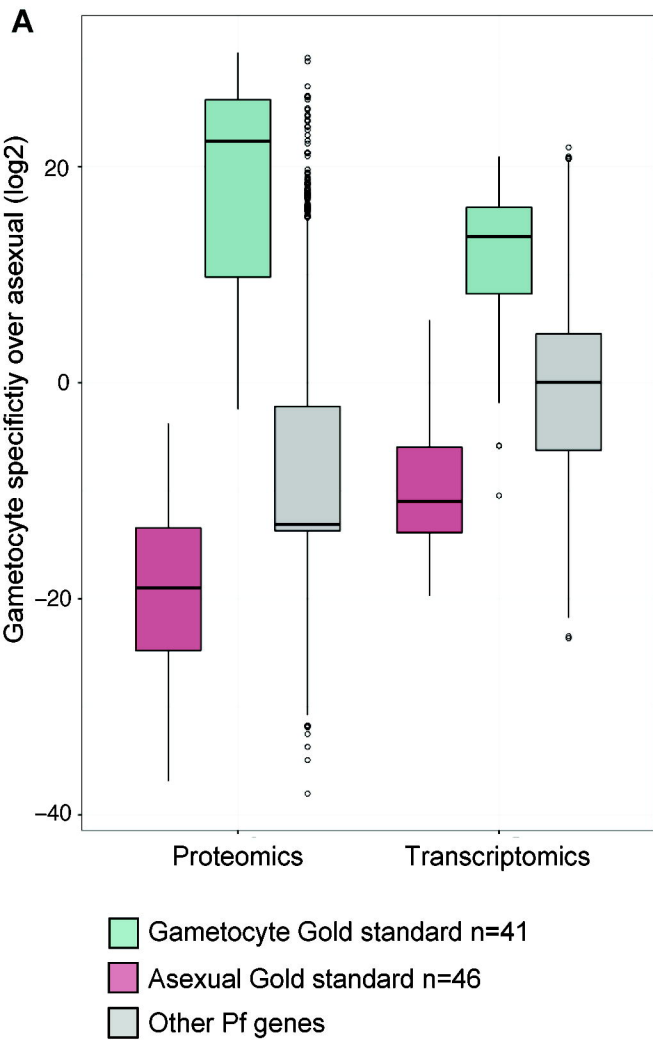
757

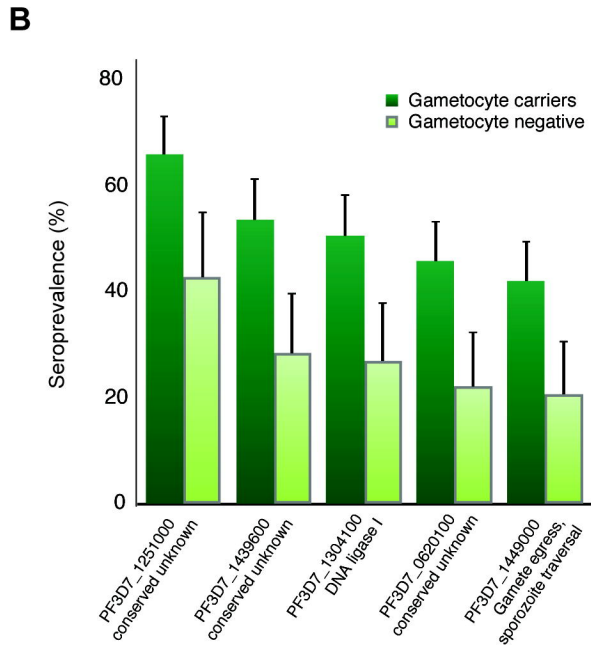
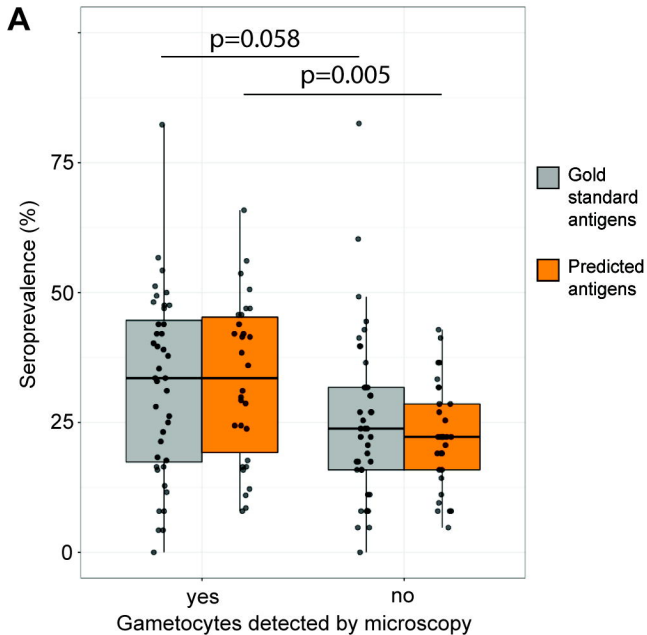


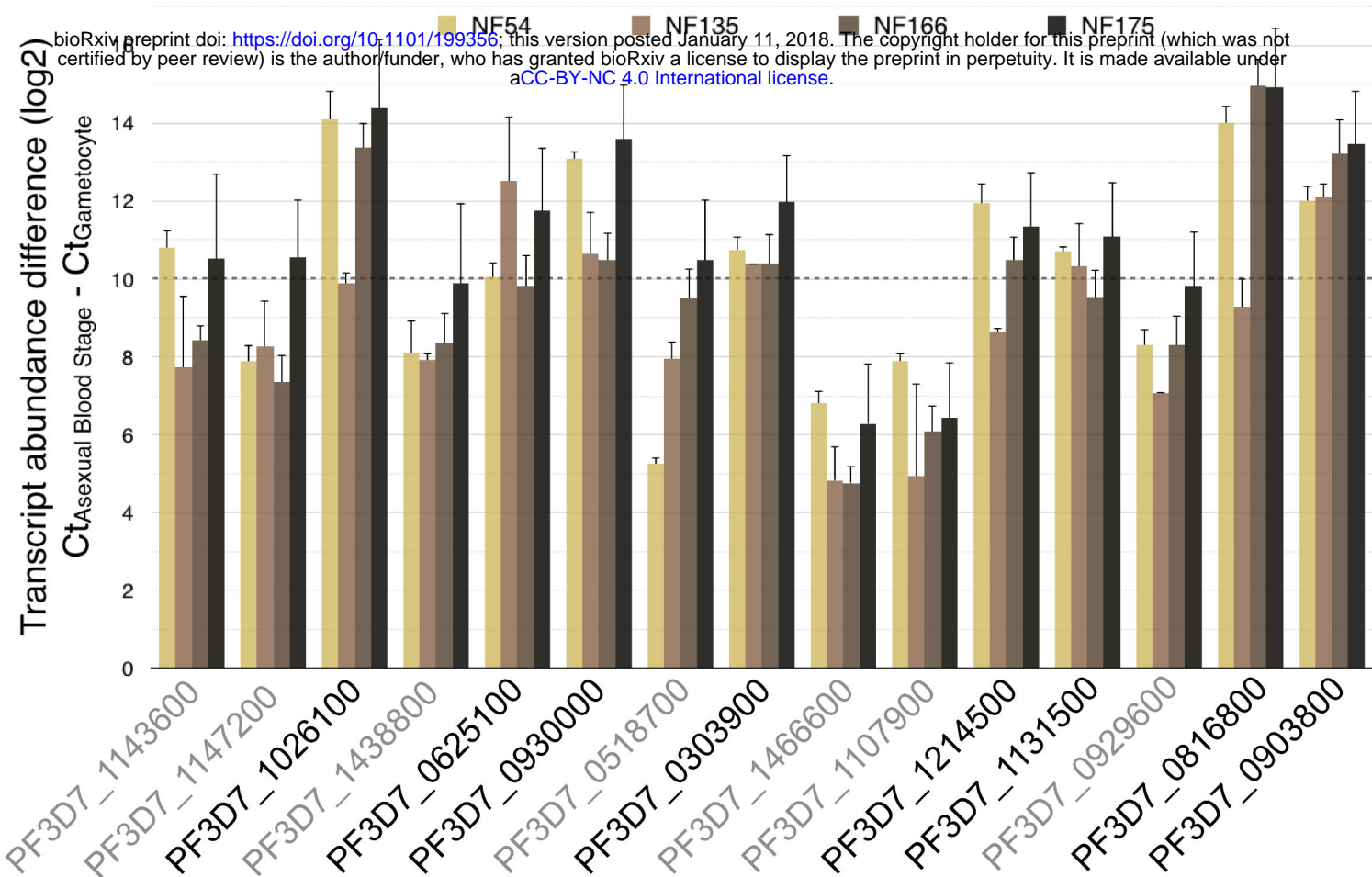
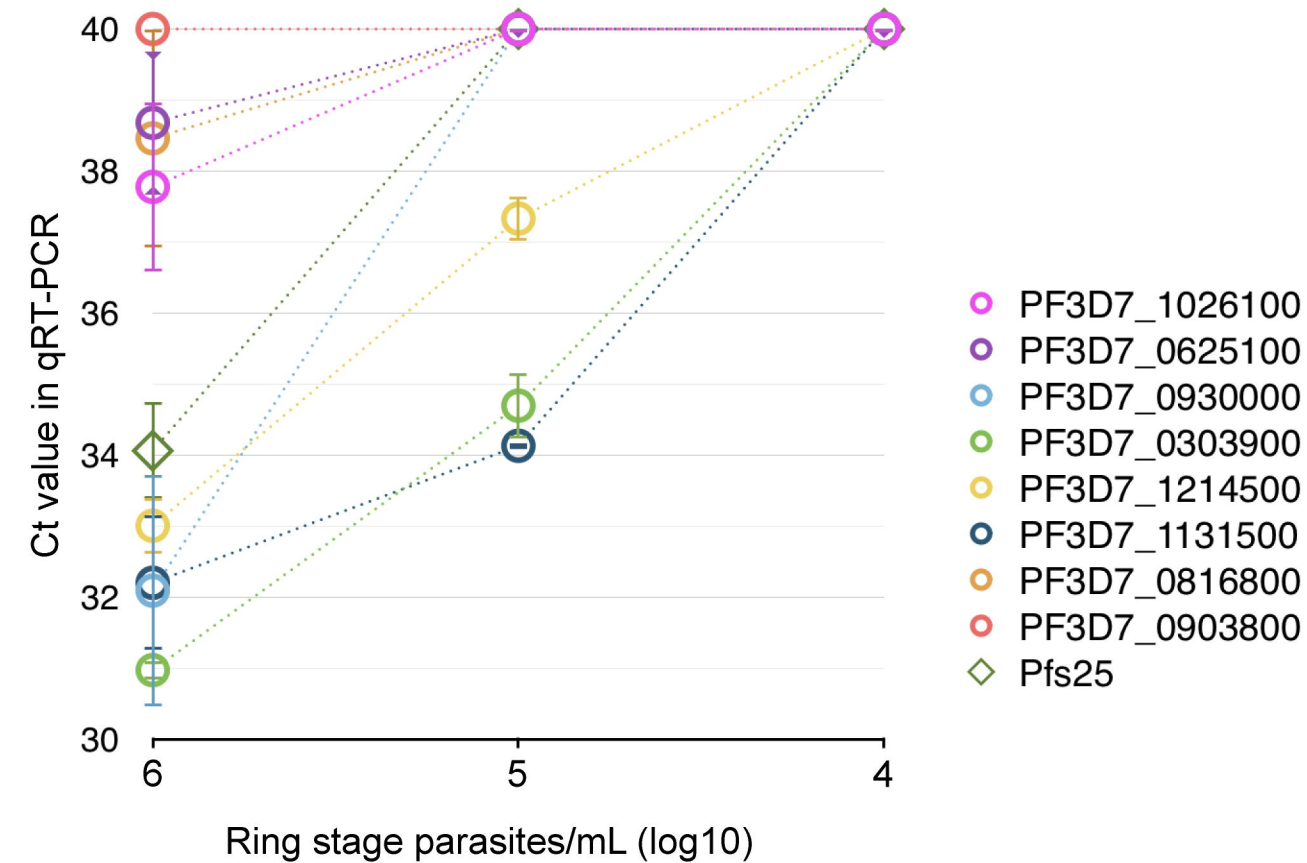
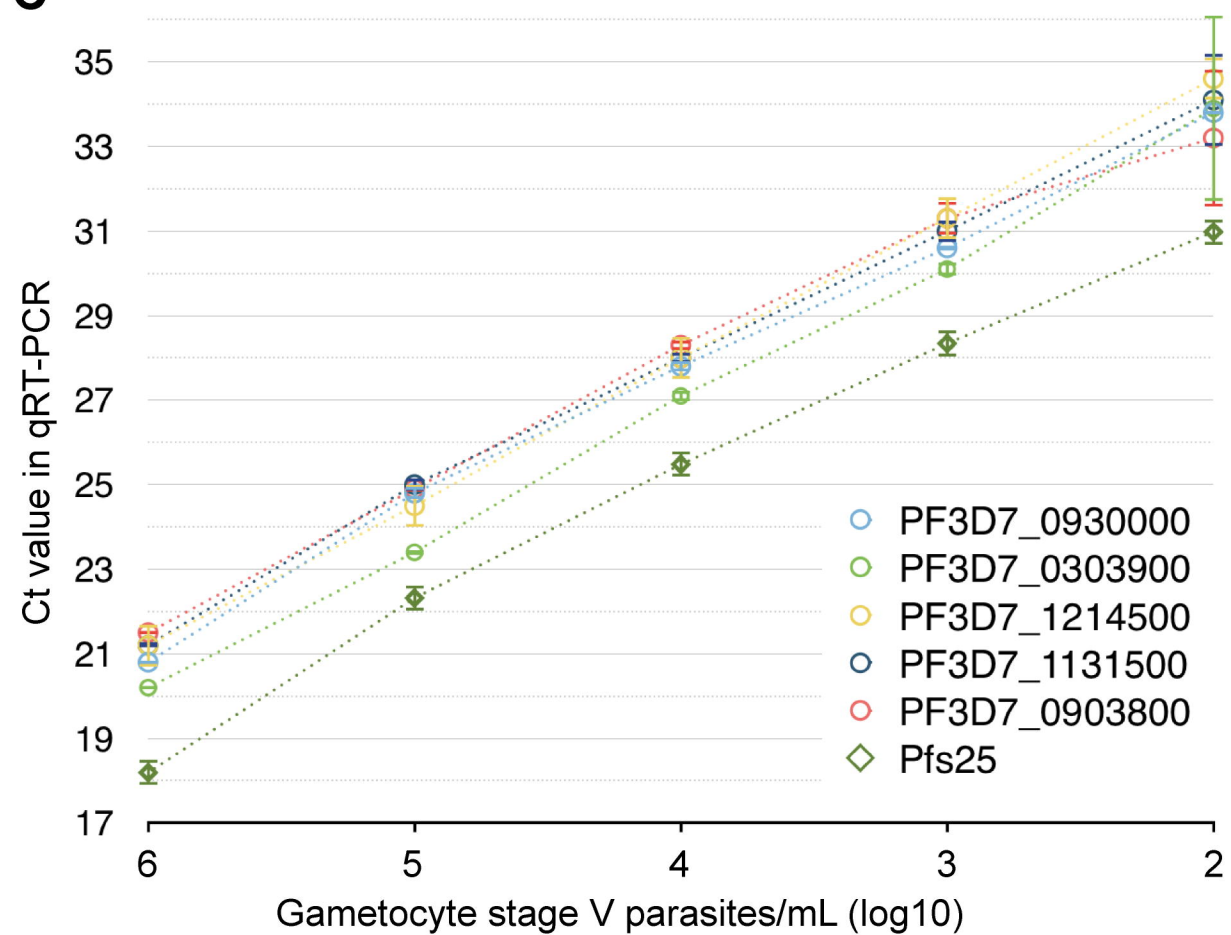
A

studies reporting as gametocyte-specific or enriched	number of proteins
0	7
1	48
2	156
3	209
4	117
5	51
6	14
total	602

B





A**B****C****D**