

1 Novel functional insights from the *Plasmodium falciparum*  
2 sporozoite-specific proteome by probabilistic integration of  
3 26 studies

4

5 Lisette Meerstein-Kessel<sup>1,2</sup>, Jeron Venhuizen<sup>1</sup>, Daniel Garza<sup>1</sup>, Emma J. Vos<sup>1#</sup>, Joshua M. Obiero<sup>3</sup>, Philip  
6 L. Felgner<sup>3</sup>, Robert W. Sauerwein<sup>2</sup>, Marynthe Peters<sup>1</sup>, Annie S.P. Yang<sup>2</sup>, Martijn A. Huynen<sup>1</sup>

7 Corresponding author: Martijn A. Huynen

---

<sup>1</sup> Center for Molecular and Biomolecular Informatics, Radboud University Medical Center, Nijmegen, The Netherlands

<sup>2</sup> Radboud Center for Infectious Diseases, Medical Microbiology, Radboud University Medical Center, Nijmegen, The Netherlands

# current address: The Hyve, Utrecht, The Netherlands

<sup>3</sup> Department of Physiology and Biophysics, School of Medicine, University of California Irvine, Irvine, CA, USA.

## 8 Abstract

9

10 *Plasmodium* species, the causative agent of malaria, have a complex life cycle involving two hosts.  
11 The sporozoite life stage is characterized by an extended phase in the mosquito salivary glands  
12 followed by free movement and rapid invasion of hepatocytes in the human host. This transmission  
13 stage has been the subject of many transcriptomics and proteomics studies and is also targeted by  
14 the most advanced malaria vaccine. We applied Bayesian data integration to determine which  
15 proteins are not only present in sporozoites but are also specific to that stage. Transcriptomic and  
16 proteomic *Plasmodium* data sets from 26 studies were weighted for how representative they are for  
17 sporozoites, based on a carefully assembled gold standard for *Plasmodium falciparum* (*Pf*) proteins  
18 known to be present or absent during the sporozoite life stage. Of 5418 *Pf* genes for which  
19 expression data were available at the RNA level or at the protein level, 1105 were identified as  
20 enriched in sporozoites and 90 specific to them. We show that *Pf* sporozoites are enriched for  
21 proteins involved in type II fatty acid synthesis in the apicoplast and GPI anchor synthesis, but  
22 otherwise appear metabolically relatively inactive, in the salivary glands of mosquitos. Newly  
23 annotated hypothetical sporozoite-specific and sporozoite-enriched proteins highlight sporozoite-  
24 specific functions. They include PF3D7\_0104100 that we identified to be homologous to the prominin  
25 family, which in human has been related to a quiescent state of cancer cells. We document high  
26 levels of genetic variability for sporozoite proteins, specifically for sporozoite-specific proteins that  
27 elicit antibodies in the human host. Nevertheless, we can identify nine relatively well-conserved  
28 sporozoite proteins that elicit antibodies and that together can serve as markers for previous  
29 exposure.

30 Our understanding of sporozoite biology benefits from identifying key pathways that are enriched  
31 during this life stage. This work can guide studies of molecular mechanisms underlying sporozoite  
32 biology and potential well-conserved targets for marker and drug development.

33

## 34 Author Summary

35

36 When a person is bitten by an infectious malaria mosquito, sporozoites are injected into the skin with  
37 mosquito saliva. These sporozoites then travel to the liver, invade hepatocytes and multiply before  
38 the onset of the symptom-causing blood stage of malaria. By integrating published data, we contrast  
39 sporozoite protein expression with other life stages to filter out the unique features of sporozoites  
40 that help us understand this stage. We used a “guideline” that we derived from the literature on  
41 individual proteins so that we knew which proteins should be present or absent at the sporozoite  
42 stage, allowing us to weigh 26 data sets for their relevance to sporozoites. Among the newly  
43 discovered sporozoite-specific genes are candidates for fatty acid synthesis while others might play a  
44 role keeping the sporozoites in an inactive state in the mosquito salivary glands. Furthermore, we  
45 show that most sporozoite-specific proteins are genetically more variable than non-sporozoite  
46 proteins. We identify a set of conserved sporozoite proteins against which antibodies can serve as  
47 markers of recent exposure to sporozoites or that can serve as vaccine candidates. Our predictions of  
48 sporozoite-specific proteins and the assignment of previously unknown functions give new insights  
49 into the biology of this life stage.

## 50 Introduction

51

52 Malaria is a mosquito transmittable disease resulting in over 220 million clinical cases and half a million  
53 deaths annually. Most deaths are caused by *Plasmodium falciparum* (*Pf*), one of the five species of  
54 *Plasmodium* that can infect humans. The infection begins with the deposition of liver-infective  
55 sporozoite forms in the skin by blood-feeding mosquitoes. These sporozoites travel to the liver where  
56 they invade, differentiate and multiply asymptotically inside hepatocytes for approximately a week  
57 before releasing red blood cell (RBC)-infective merozoites into the circulation. The subsequent asexual  
58 multiplication, rupture and re-invasion of the parasites into circulating RBCs cause the symptoms  
59 associated with malaria.

60 Identifying specific sporozoite proteins will aid in the understanding the biology of this highly motile  
61 stage in which the parasite is directly exposed to the body's immune system, similar to blood stage-  
62 infective merozoites. Sporozoites deposited in the skin can take up to 90 minutes to reach the liver [1],  
63 a much longer potential exposure to immune system than the 90 seconds of merozoites. Mosquito  
64 midgut-derived and circulating sporozoites are less infectious for hepatocytes than those that have  
65 resided in the salivary gland [2, 3]. The ability to improve and remain infectivity in the salivary gland of  
66 mosquitoes for an extended period of time (approximately 1 week) is an intriguing phenomenon that  
67 is poorly understood [4]. It is suggestive of a molecular landscape where the parasite is kept in low  
68 activity but can quickly activated to evade immune systems, invade and develop in hepatocytes.  
69 Understanding the specific molecular make-up of sporozoites will shed light on this aspect of its biology  
70 and may reveal new targets for interventions.

71 A number of studies have identified genes that are expressed in sporozoites during both their  
72 development in mosquito midgut (oocyst) and in the salivary gland. They have identified expression at  
73 the RNA level[5] and at the protein level[6] with a specific emphasis on surface proteins[7]. The studies  
74 varied in their focus and resolution, with RNA level studies facing the challenge of their relevance for  
75 protein expression[8] while proteomics studies face challenges in detecting low abundance proteins.  
76 More importantly, while such studies address the question what is present in the sporozoite, they do  
77 not address what is specific to it. Therefore, to prioritize sporozoite specific proteins, we implemented  
78 a naïve Bayesian data integration in which both transcriptomics and proteomic data were included.  
79 For informed data integration, gold-standard lists representing proteins that are either sporozoite  
80 specific or non-sporozoite specific were manually assembled from the existing literature on individual  
81 proteins. Using this gold standard, published *Plasmodium* datasets were weighted by the presence of  
82 sporozoite specific proteins and the absence of proteins known to be absent from sporozoites,  
83 allowing us to obtain an extended list of predicted sporozoite specific proteins. We classified 90

84 proteins as sporozoite-specific, of which 67 were not part of the gold standard. “Conserved,  
85 hypothetical proteins with unknown functions” were examined using sensitive homology and  
86 orthology detection tools [9, 10] to predict their function and shed light on the biology of sporozoites.

87 From the proteins that in the assembled proteomics data were identified to be present in sporozoites,  
88 we examined whether we can identify a limited set of proteins for which human antibodies are  
89 generated[11] and that can serve as or markers of exposure to sporozoites. As attractive targets would  
90 be conserved between different *P. falciparum* strains, we sequenced three *P. falciparum* strains, NF54,  
91 NF135 and NF166, and required selected proteins to have limited genetic diversity between those  
92 strains as well as in sequenced genomes in PlasmoDB[12].

93

## 94 Methods

95

### 96 **Sporozoite protein and transcriptomic data sets**

97 The data integration was performed using transcriptomic and proteomic data sets from 22 studies  
98 describing all life cycle stages of *P. falciparum* (**S1 Table**). Data sets were obtained from PlasmoDB  
99 version 43 [12], and literature (**S1 Table**). In addition to that, transcriptomics data on the liver and  
100 blood stage of *Plasmodium cynomolgi* and the sporozoites of *Plasmodium vivax* were included, as  
101 well as proteomics studies covering the sporozoite stage of *P. vivax* and the liver stage of *P. yoelii*,  
102 respectively, resulting in a total of 48 data samples from 26 different studies (**S1 Table**). Data from  
103 non-*falciparum* studies were converted into *P. falciparum* IDs using the orthologs lists available on  
104 PlasmoDB, favouring orthologs that are syntenic in the case of multiple options.

105

### 106 **S1 Table: Overview of all datasets used for the Bayesian data integration**

107

### 108 **Gold standards**

109 Bayesian data integration requires gold standards, in this case of proteins known to be highly enriched  
110 in sporozoites or depleted from them. We required positive gold standard proteins to be present  
111 dominantly in sporozoites based on western blot or immunofluorescent assay data, resulting in a  
112 selection of 31 sporozoite-specific proteins (**S2 Table**).

### 113 **S2 Table: Positive gold standard members**

114 Evidence for expression in sporozoites

115

116 The negative gold standard was curated by searching literature for non-sporozoite proteins. We  
117 selected 19 proteins based on their literary evidence of absence from sporozoites. Additionally, 20  
118 gametocyte specific proteins [13] were added to the negative gold standard, increasing the total  
119 number to 39 negative gold standard proteins (**S3 Table**). These proteins spanned the remaining *P.*  
120 *falciparum* life cycle stages in the human host, with the majority being found in (a)sexual blood  
121 stages.

122

### 123 **S3 Table: Negative gold standard members**

124 Evidence for expression in other life stages

125

## 126 **Bayesian data integration**

127 The used data sets were examined for their correlations with each other (**S1 Figure**). By and large  
128 most data sets show little correlation. We did leave the few correlated data sets in to keep the data  
129 integration transparent and maximize the amount of included information. Oocyst-derived and  
130 salivary gland sporozoites showed high correlations with each other, which led us to combining all  
131 respective studies into the “sporozoite” data input and not make a distinction between those stages.

### 132 **S1 Figure: Correlations of all integrated data sets with each other, hierarchically clustered.**

133 Study code (two letters and year) as in table S1. Samples are coded as sporozoite (sporo), salivary  
134 gland sporozoite (SGS), oocyst derived sporozoite (ODS) or oocyst (OOC), liver stage (LS), blood stage  
135 (BS, ery, asexual, merozoite, schiz, ring) or other stages (gametocyte/gam, zygote, ookinete).  
136 Transcriptomics studies are given a percentile (percent) for each gene they detected and proteomics  
137 studies a score for each unique peptide (uniq\_pept) per gene.

138

139 Proteomic data were converted into unique peptide counts for each protein identified and  
140 transcriptomic data were converted into expression percentiles for a total of 5668 *P. falciparum* gene  
141 IDs. Proteomic and transcriptomic data was binned consistently for all data sets, with 0, 1, or >1  
142 identified unique peptides, or into four bins, containing transcripts that are in the >80 percentile, >60  
143 percentile, >40 percentile and <40 percentile, respectively. The data sets were then weighted  
144 according to their ability to retrieve the gold standard proteins. Each bin in each data set was given a  
145 log<sub>2</sub> score according to equation (1), where B = present in bin, S = sporozoite specific and nonS = not  
146 sporozoite specific.

$$147 \log_2(P(S|B)) = \log_2\left(\frac{P(B|S)}{P(B|nonS)}\right) \quad (1)$$

148 The Bayesian score for an individual protein is then the sum of the scores for the bins in which it  
149 occurs (one bin per data set). As the expected number of sporozoite-specific proteins was unknown,  
150 no *prior* was included, rather we used cutoffs based on the position of known sporozoite specific  
151 proteins to define sporozoite specific proteins and sporozoite enriched proteins.

### 152 **Overrepresentation of function categories in sporozoite proteins**

153 GO terms were acquired from PlasmoDB and formatted into a .gmt file according to the format  
154 specified by the GSEA server at the Broad Institute[14]. GSEAPreranked on the ranked list of  
155 sporozoite-specific proteins was used with the conservative “preranked” option in the “classic mode”,

156 i.e. using Kolmogorov Smirnov statistics to determine enriched GO terms. The large variable protein  
157 families RIFIN, STEVOR and PfEMP1 (*var* genes) were left out of the analysis. A new GO term for gliding  
158 motility in *Plasmodium* was assembled by searching literature for proteins associated with gliding  
159 motility in sporozoites, ookinetes and merozoites. For this we extended the list of glideosome  
160 associated proteins assembled by Lindner *et al.*[15] with 10 new proteins possibly associated with the  
161 glideosome. These proteins were either annotated with the motility GO term GO:0071976 (SSP3,  
162 CelTOS, LIMP protein, plasmepsin VII and glideosome-associated connector), or otherwise known to  
163 be involved in gliding motility (CTRP[16], SIAP1[17], GAPDH[18], GAP40[19], IMC1I[20]). Similarly we  
164 added a list of GPI-anchored proteins to the set of processes for which we examined enrichment using  
165 the list of Gilson *et al.*[21], based on the hypothesis that the type II fatty acid synthesis was required  
166 for the creation of GPI anchors [22]. To investigate Pfam domain enrichment in sporozoite proteins,  
167 Pfam annotations for all proteins were downloaded from PlasmoDB and were used as gene sets in the  
168 GSEA. To prevent the GSEA analysis results for enriched pathways to be affected by the large numbers  
169 of hypothetical proteins in *P. falciparum*, those were filtered out before the analysis. Furthermore,  
170 proteins that were part of the gold standard were left out of the analysis to prevent circular arguments.  
171 Proteins for which no transcriptomic and proteomic data were available were left out as well. Manual  
172 examination of the Gene Set Enrichment Analysis results showed that some gene sets were  
173 significantly enriched in sporozoites only because they were depleted from the low scoring proteins.  
174 Therefore, we required for the significantly enriched processes in sporozoites that they actually  
175 contained proteins in the set of “sporozoite enriched proteins.”

176

### 177 **Sequencing of *P. falciparum* strains**

178 NF54 originates from West Africa, and is isolated from a woman infected in the Netherlands nearby  
179 an airfield[23]. The NF135 clone is a clinical isolate that originates from Cambodia[24]. The NF166  
180 clone is a clinical isolate from a child that visited Guinea[25]. Whole genome sequencing of the three  
181 strains was performed with Illumina NextSeq 500, resulting in raw paired-end fastq reads of 151 base  
182 pairs (bp).

### 183 **Quality control and trimming**

184 To examine the quality of the raw fastq reads, FastQC (version 0.11.5) was used[26]. The Nextera  
185 Transposase sequence contamination at the 3' ends of the reads were trimmed off with a stringency  
186 of 12 bp, using Trim Galore (version 0.4.3)[27]. CleanNextSeq\_paired was used to remove excess of  
187 G's from 3' ends of the reads after 100 bp[28]. Only reads with a minimal length of 35 base pairs  
188 were retained.



## 189 **Alignment and variant calling**

190 The *P. falciparum* 3D7 reference genome (v3.0, PlasmoDB, plasmodb.org, 14 chromosomes,  
191 mitochondrial genome and apicoplast genome) was indexed and the trimmed fastq reads were  
192 aligned to the reference genome using Bowtie2 (version 2.2.8) with the local alignment setting[29].  
193 The SAM files obtained were converted to BAM files and subsequently sorted and indexed with  
194 SAMtools (version 1.4.1)[30]. SNPs and indels were called using SAMtools and BCFtools (version  
195 1.4.1). Alignments were visually inspected with the Integrative Genomics Viewer (IGV, version  
196 2.3.98)[30]. Coverage was calculated with BEDtools (version 2.26.0)[31].

## 197 **Filtering of SNPs and indels**

198 Filtering of SNPs and indels was performed with BCFtools (version 1.4.1)[30]. SNPs with a base Phred  
199 quality (Q) > 30 were used for further analysis. Furthermore, we required that the proportion of high  
200 quality bases (the DP4 scores in the VCF files) supporting the indel or the SNP in  $\geq 75\%$  of the called  
201 bases to include them for further analysis.

## 202 **Effect of SNPs and indels on protein sequences**

203 The effect of the mutations on the predicted protein sequences was determined using the Genomic  
204 Ranges package in R [32]. Amino acid sequences of all proteins were compared to amino acid  
205 sequences of the reference genome and with each other. Any variation at the amino acid level  
206 between reference and the examined strain and among the examined strains was denoted.

207

## 208 **Homology detection**

209 Homology detection of *P. falciparum* proteins with unknown function was done using HHPred[9] with  
210 default settings (HHblits, 3 iterations). For orthology detection we used best bidirectional best hits at  
211 the level of sequence profiles[10].

## 212 **Detecting a set of proteins that elicit antibodies in all CPS volunteers**

213 We used a greedy “set cover” algorithm that in each step selects, from a list of evolutionary  
214 conserved sporozoite proteins (less than eight nonsynonymous SNPs per kb in PlasmoDB), the one  
215 that is immunogenic (elicits antibodies) in the highest numbers of volunteers for which no  
216 immunogenic protein was already in the set. In the analysis from Obiero *et al.* [11], some long  
217 proteins were split into separate peptides. Those were analyzed separately by the algorithm: i.e. as if  
218 they were separate proteins. When multiple proteins were immunogenic in the same number of  
219 volunteers for which no immunogenic protein had been selected yet, we chose from those the one

220 that had the highest immunogenicity in all the volunteers. We furthermore required genes to have  
221 maximally two non-synonymous SNPs in the combined NF135 and NF166 strains.

222

223

## 224 Results

225

### 226 **Data integration**

227 At least 26 studies, containing 48 data sets have been published of *P. falciparum*, *P. cynomolgi*, *P. vivax*  
228 and *P. yoelii* transcripts and proteins at various developmental stages, including 11 that were specific  
229 to the sporozoite stage (**S1 Table**). In order to optimally exploit those data to obtain sporozoite-  
230 enriched proteins we integrated them in a Bayesian manner (Methods). Integrating the data sets using  
231 the sets of 31 positive and 39 negative gold standard proteins (**S2 Table, S3 Table**) produced a list of  
232 all proteins in *P. falciparum* ranked according to their likelihood of being sporozoite-specific (**S4 Table**).  
233 The score distribution of the negative and positive gold standard proteins varied depending on using  
234 all available data sets, or proteomic or transcriptomic data separately (Figure 1A and **S2 Figure**). The  
235 Bayesian integration using only transcriptomic data sets resulted in a ranked list where 14 negative  
236 gold standard genes scored higher than the lowest scoring positive gold standard gene (**S2 Figure**). This  
237 overlap was lower when using only proteomic data sets (**S2 Figure**), but still contained 8 proteins. The  
238 least overlap between positive and negative gold standard members was observed when combining  
239 transcriptomic and proteomic data (Figure 1A). This also produced an outlier, STARP. Although this  
240 protein was identified as being a sporozoite protein [33], in our combined data sets it scored 42 times  
241 lower than the second to last scoring positive gold standard protein. It therewith does not appear  
242 specific to sporozoites in the available data. STARP was excluded from the score distribution of gold  
243 standard proteins that was used to define the sporozoite specific proteins (see below), but gold  
244 standard list and integration where not changed *post hoc*. Based on the observed overlaps, we decided  
245 to continue our research using the ranked list based on the combined proteomic and transcriptomic  
246 data sets.

247

### 248 **S2 Figure: Distribution of sporozoite-specific probability score.**

249 Score for positive and negative gold standard and all remaining Plasmodium falciparum  
250 genes (rest) when using only proteomics data as input (A) or only transcriptomics data (B).

251

252 **S4 Table: Ranked genes for their sporozoite-specificity**

253

254 **Sporozoite specific proteins and sporozoite enriched proteins**

255 Proteins were considered *sporozoite-specific* when ranking in the first quartile of the gold standard  
256 protein list i.e. scoring above 12.27, which represents a factor of being at least  $2^{12.37}$  to  $2^{5200}$  more  
257 likely to be specific to sporozoites than to any of the other stages. Do note that hereby we ignore the  
258 *prior* probability of a protein being sporozoite specific at all. Although such a prior is standard in  
259 Bayesian data integration, in this case a prior, e.g. of 1/5, is given the high cut-off for sporozoite specific  
260 proteins that we used, not very relevant. As we did not consider STARP to be sporozoite specific (see  
261 above), we considered proteins that scored higher than the second to lowest positive gold standard  
262 protein (LIMP protein, score = -1.31), to be *enriched* in sporozoites. Finally, the abundance of unique  
263 peptides by mass-spectrometry was assessed for each protein. Proteins were deemed *present* in  
264 sporozoites when identified in two independent studies or with more than 1 unique peptide in at least  
265 one study. Our analysis thus identified 90 sporozoite-specific proteins, 1105 sporozoite-enriched  
266 proteins and 2736 that were present in sporozoites (**S4 Table**). Out of the 90 sporozoite-specific  
267 proteins, 67 were not part of the positive gold standard list. We validated our predictions by 5-fold  
268 cross validation (5-fold CV) by randomly skipping 1/5<sup>th</sup> of the gold standard proteins from the data  
269 integration and assessing their predicted sporozoite specificity based on the remaining data (Fig 1B).  
270 The high sensitivity and specificity indicated that novel sporozoite-specific proteins would also score  
271 higher than non-sporozoite specific proteins. We also compared the ranking of the gold standard  
272 proteins based on the integrated data with a ranking based on individual data sets of sporozoite RNAs  
273 and proteins. The cross validation separated the gold standard proteins better than individual data  
274 sets, supporting the integration of multiple data sets (Fig 1B).

275

276 **Fig 1. Bayesian data integration identifies sporozoite-specific genes in *P. falciparum***

277 **A)** Bayesian score distributions of the proteins from the negative gold standard, the positive gold standard and  
278 the remaining proteins, **B)** cross validation of our predictions (5-fold cross validation) with all data sets (Black  
279 solid line) and predictions with individual data sets (dashed coloured lines).

280

281 **Function prediction of non-annotated proteins**

282 Many of the genes in the *P. falciparum* genome encode hypothetical proteins with unknown molecular  
283 function [34]. The fraction of unknowns that is specific for sporozoites (32 out of 90, 36%) is at least as  
284 high as in the rest of the genome (32%). To improve understanding of their potential functions, we  
285 examined the proteins for domains with known functions from any species using sensitive homology  
286 detection with HHpred[9], combined with manual examination of conserved residues. We compared  
287 sporozoite specific proteins with PFAM domains, the human proteome, and the proteome of  
288 *Toxoplasma gondii*, an apicomplexan related to *P. falciparum* that is present in the HHpred database.  
289 The sporozoite specific protein with the highest score that was not part of the gold standard,  
290 PF3D7\_0104100, showed (barely) significant sequence similarity with the prominin family ( $E=10E-4$ )  
291 and low levels of sequence identity with e.g. human prominin-2 (12%). To cross-check the homology  
292 of PF3D7\_0104100 with prominin we examined homology with *T. gondii* proteins. The sequence  
293 similarity between the *P. falciparum* protein and *T. gondii* TGME49\_218910 ( $E=3.4e-44$ ) and between  
294 the *T.gondii* protein and the human prominin-2 were highly significant ( $E=2.3e-21$ ), indicating that  
295 PF3D7\_0104100 is indeed member of the prominin family. PF3D7\_0104100 has like the prominin  
296 family five (predicted) transmembrane regions with most of the protein localized outside the cell  
297 (Figure 2). Analysis of a sequence alignment with orthologs in other *Plasmodium* species reveals the  
298 conservation of ten cysteine residues that are all predicted to be extracellular (Figure 2). Such  
299 extracellular cysteines can form disulphide bonds as has been observed for other extracellular  
300 *Plasmodium* protein domains [35] and has been suggested for human prominin[36]. Two pairs of the  
301 extracellular cysteines were conserved between the human prominin and PF3D7\_0104100 (Figure 2,  
302 **S3 Figure**). We were able to predict molecular functions of three other sporozoite specific proteins and  
303 27 sporozoite enriched proteins using HHpred and best bidirectional hits with human proteins at the  
304 level of sequence profiles (**S6 Table**)[10].

**Fig 2 Predicted membrane topology of Pf3D7\_0104100, a sporozoite-specific protein that is homologous to the prominin/CD133 protein family.**

The level of polymorphisms among *P. falciparum* strains is indicated for the separate regions as the average number of polymorphisms per nucleotide in the strains in PlasmoDB. PF3D7\_0104100 has a high density of polymorphisms within *P. falciparum* strains that are concentrated in the second extracellular loop. Cysteines that are conserved among the homologs in *Plasmodium* species are indicated. The cysteines that are conserved also in human homologs are in bold. Note that the conserved cysteines occur in close proximity to each other, suggesting the formation of disulphide bonds.

305

306 **S3 Figure: Alignment of PF3D7\_0104100 from *P.falciparum*, TGME49\_218910 from *Toxoplasma***  
307 ***gondii* and Prominin-2 from *Homo sapiens*.** The in PF3D7\_0104100 predicted transmembrane  
308 regions are underlined, the conserved cysteines are boxed. Do note that the fourth transmembrane  
309 region is relatively long. Predictions with other tools than TMHMM [1] like Phobius [2] indicate a

310 shorter TM region, putting the cysteine that is located in that TM region in the extracellular space.  
311 The *Toxoplasma* protein was included because its sequence profile has significant sequence similarity  
312 against both the human protein profile ( $E= 2e-20$ ) and the *Plasmodium* protein profile ( $E= 3.4e-44$ ),  
313 while the similarity between the human and the *Plasmodium* protein is less significant ( $E=0.0001$ ).

314

315 **S6 Table: Sporozoite enriched proteins annotated as “unknown function” in PlasmoDB for which**  
316 **orthology with human proteins could be detected using best bidirectional hits at the level of**  
317 **sequence profiles**

318

319

### 320 **Over-represented domains and pathways**

321 We examined whether specific protein domains and pathways were over- or under-represented in the  
322 sporozoite proteins using Gene Set Enrichment Analysis[14]. For this we augmented the list of proteins  
323 involved in specific processes from PlasmoDB with the glideosome and glycosylphosphatidylinositol  
324 (GPI) anchor proteins that we considered specifically relevant to sporozoites (Methods)[14]. At the  
325 domain levels there was an over-representation of three protein domains: PF09175 (also named  
326 Plasmod\_dom\_1, unknown function), PF08373 (the RAP domain, putatively RNA binding) and PF12879  
327 (the C-terminal domain of the SICA proteins that are associated with parasitic virulence). Two of these  
328 three domains are either unique to *Plasmodium* species (PF12879) or to *P. falciparum* (PF09715). At  
329 the level of pathways, we observe significant enrichment of the Type-II fatty acid synthesis (FAS-II) and  
330 the GPI anchor biosynthesis pathways. The enzymes present in the FAS-II pathway are located in the  
331 apicoplast – a reduced plastid-like organelle that shares similarity with chloroplasts of algae and plants  
332 [37]. This over-representation is consistent with its essentiality for sporozoite development [38].  
333 Interestingly, *P. falciparum* has fourteen Acyl CoA synthetases [39], of which the distribution over the  
334 various fatty acid metabolizing processes is largely unresolved. There is one Acyl CoA synthetase, ACS2,  
335 that is specific to sporozoites and that would be an interesting candidate to convert the fatty acids  
336 produced in the apicoplast to acyl-CoA. Curiously, there was no enrichment in sporozoites of members  
337 of the acyl-ACP thioesterase family (PF3D7\_1135400 and PF3D7\_0217900). Acyl-ACP thioesterases  
338 play an important step in the pathway by hydrolyzing the acyl moiety from the ACP before it can be  
339 converted to acyl-CoA.

340 The second enriched pathway is GPI-anchor biosynthesis, with three proteins (GPI  
341 mannosyltransferases 1, 2 and 3: PF3D7\_1210900, PF1247300 and PF3D7\_1341600) enriched in  
342 sporozoites. GPI represents a class of glycolipids found as either free lipids or attached to proteins [40,  
343 41]. Surface parasite proteins (such as the merozoite surface proteins 1, 2 and 4) anchor to the parasite

344 cell-membrane via GPI moieties [42]. Sporozoites shed surface proteins (such as circumsporozoite  
345 protein and thrombospondin related anonymous protein) when moving [43]. With GPI being used to  
346 anchor proteins to the parasite surface, it would have to be synthesized constantly to replace the  
347 surface proteins shed during motility. There was no enrichment of the glideosome or of GPI anchor  
348 proteins among proteins specific to the sporozoite. However, 39 of the 52 glideosome associated  
349 proteins are “present” in the sporozoite stage (75%), representing a slight overrepresentation when  
350 compared to the 2,736 of 5,447 proteins (50,2%) detected in sporozoites (Fisher’s exact, P=0.0018).  
351 Similarly, there are 19 out of 28 GPI-anchor proteins identified in Gilson *et al.* [21] among the proteins  
352 present at the sporozoite stage, but they are not specifically enriched in sporozoites and hence are  
353 also expressed in other developmental stages. The fact that proteins involved in movement (such as  
354 glideosome and GPI) are enriched but not specific to sporozoites indicated that sporozoite share a  
355 common movement machinery with other the life-cycle stages.

### 356 **Protein functions of sporozoite-specific proteins in the context of sporozoite biology**

357 Apart from examining over represented pathways we also examined the individual sporozoite specific  
358 proteins with regard to their potential function in sporozoite biology. One of the interesting features  
359 of sporozoites is their inactivity in the sporozoite glands, which is among others maintained by  
360 translational repression with the Puf2 protein [44]. In view of the observation that the translational  
361 repression occurs in two waves [44], it is interesting that the sporozoite specific protein  
362 PF3D7\_0411400 contains a Dead box helicase like its paralog of DOZI that is involved in translational  
363 repression in gametocytes of *Plasmodium berghei* [45]. From a gene expression regulatory perspective  
364 are furthermore interesting: three Zinc finger proteins (PF3D7\_0615600, PF3D7\_0521300,  
365 PF3D7\_1008600), and the transcription factor AP2-04 that has mainly been implicated in ookinetes  
366 [46]. With respect to the epigenetic regulation, the histone deacetylase sir2b (PF3D7\_1451400) that is  
367 involved in epigenetic silencing of var gene expression [47] is sporozoite specific, in contrast to its  
368 paralog sir2a that is not enriched in sporozoites. Examining the sporozoite specific signaling proteins  
369 with respect to their potential function in maintaining the temporal inactivity of the parasite we noted  
370 the RAC-beta serine/threonine protein kinase PfAKT (PF3D7\_1246900) and two 14-3-3 domain  
371 proteins (PF3D7\_1422900 and PF3D7\_1362100). The PfAKT that is involved in Artemisinin resistance  
372 [48] is specifically interesting in view of the role of its metazoan orthologs, AKT1/2/3, in regulating  
373 cellular metabolism to support cell survival [49]. In human it phosphorylates a large number of targets,  
374 a number of which are after phosphorylation bound by 14-3-3 proteins. Two 14-3-3 proteins are  
375 sporozoite specific: Pf14-3-3I (PF3D7\_1422900) and PF3D7\_1362100. Pf14-3-3I has been shown to  
376 bind phosphorylated Histone H3 [50]. The Raf kinase inhibitor RKIP (PF3D7\_1219700), which may be a  
377 substrate of the calcium -dependent protein kinase 1 [51], may also be involved in keeping sporozoites

378 in an inactive state as its human ortholog in inhibits the MAPK pathway [52]. Furthermore, with  
379 respect to inactivity of sporozoites it is interesting to mention that Prominin-1/CD133, one of the two  
380 human orthologs of the Prominin gene (PF3D7\_0104100), is known as a marker for dormant stem cells,  
381 e.g. in melanoma [53] or in Neural cells [54], and has been shown to activate the PI3K/AKT  
382 pathway[55]. Together, these data suggest that these proteins may also modulate sporozoite survival  
383 via a mammalian Akt-like pathway as previously shown in asexual blood stages.  
384 Finally, the presence of the tubulin polymerization promoting protein p25 alpha (PF3D7\_1236600) as  
385 well as thioredoxin-like protein (PF3D7\_0919300) and thioredoxin-like associated protein 1  
386 (PF3D7\_1230100) whose orthologs are associated with tubulin in *T. gondii* [56] are worth noting in  
387 view of the essential role of the microtubules in the sporozoite stage [57]. In contrast to these  
388 proteins, the  $\beta$ -tubulin, the two  $\alpha$ -tubulins and the other thioredoxin-like associated proteins  
389 associated with tubulin [56] are not enriched in sporozoites. Given the presence of  $\beta$ - and  $\alpha$ -tubulins  
390 in all stages of the life-cycle, PF3D7\_1236600, PF3D7\_0919300 and PF3D7\_1230100 may be crucial  
391 players involved in providing sporozoites with their characteristic, thin cylindrical shape.

392

### 393 **Fig 3. Apicoplast fatty acid synthesis proteins are enriched among sporozoites.**

394 The width of the arrows is determined by the Bayesian score reflecting the level of over representation  
395 of that enzyme in sporozoites, e.g. 16 for ACS2 and 8 for FabB/F (**S4 Table**). For PDH that consists of  
396 three proteins, the width of the arrow was determined by the average of those three. Most of FASII  
397 proteins are enriched in sporozoites, except PKII, FABZ and LipA. The scheme is a simplification of the  
398 pathway as depicted by Shears *et al.* [22], to which ACS2 was added as it is highly enriched in sporozoites  
399 and relevant for fatty acid synthesis.

400

### 401 **Under-represented pathways**

402 In contrast to the paucity of upregulated processes, there is significant under representation (FDR  $\leq$   
403 0.001) of processes linked to splicing, translation, translation elongation, folding of proteins as well as  
404 proteolysis (S5 Table). These processes are primarily modulated by a number of sporozoite specific  
405 proteins involved in transcriptional silencing e.g. PF3D7\_0411400 that contains a Dead box helicase  
406 domain[45] and PF3D7\_1451400, a histone deacetylase involved in the epigenetic silencing of  
407 virulence (*var*) gene expression [47]. Furthermore, there is significant depletion of proteins involved in  
408 carbon metabolism: glycolysis and citric acid cycle (S5 Table).

409

### 410 **S5 Table: Biological processes enriched (top list) or depleted (bottom list) in the sporozoite**



411 Enrichment based on the relative absence of the proteins involved at an FDR  $\leq 0.01$  and an  
 412 enrichment score  $> 0.20$ . The type II Fatty Acid Synthesis are the genes from Figure 2, of which the  
 413 gold standard genes were left out in the GSEA analysis. The “translocation of peptides or proteins  
 414 into hosts” GO category did not have any proteins among the sporozoite enriched proteins, and was  
 415 left out of the description of the results. (N)ES: (normalized) enrichment score, FDR: false discovery  
 416 rate

417

#### 418 **Selecting sporozoite proteins as targets or markers of past infection**

419 Antibodies that bind sporozoite proteins have been induced in 38 volunteers after chloroquine  
 420 chemoprophylaxis with *P. falciparum* sporozoites (CPS)[11], Table S7. Induced antibody profiles  
 421 represent a blue print of immunogenic proteins in sporozoites, liver stages and early blood stages.  
 422 Here we focus on the set of sporozoite target proteins that may be used as markers of previous  
 423 sporozoite exposure or may act as potential targets for vaccines. Minimal sequence variation  
 424 between *Pf* strains would thereby strengthen the candidature of the proteins for epidemiological or  
 425 clinical applications, however antibody eliciting proteins including CSP, show relatively high levels of  
 426 polymorphisms among sequenced malaria strains (**S4 Figure**), and also sporozoite specific proteins  
 427 show high levels of polymorphisms (**S5 Figure**). The selection of proteins that could serve as markers  
 428 of exposure or vaccine candidates is a compromise between on the one hand protein sequence  
 429 conservation and on the other hand the frequency of volunteers with antibodies to that protein. As  
 430 the maximum level of sequence variation between candidate marker proteins we used 8 non-  
 431 synonymous SNPs per kilobase, which is lower than the variation in for instance Pfs48/45, a highly  
 432 conserved gametocyte protein with 8.9 nonsynonymous SNP per kilobase. As the minimum number  
 433 of people in which a protein should elicit antibodies we chose six (out of the 38). Using a “greedy  
 434 search algorithm” (methods) we selected a set of nine proteins of which at least one elicited  
 435 antibodies per volunteer (Table 1).

#### 436 **Table 1: putative markers of exposure to sporozoites**

437 Genes selected by the greedy method to cover all volunteers with their gene ID, function and  
 438 number of volunteers with antibodies after CPS immunization as reported by Obiero *et al.* [11]. The  
 439 non-synonymous SNPs are given as a proxy for genetic variability, with PlasmoDB and two sequenced  
 440 laboratory strains as reference.

Gene ID	function	No. (%) of volunteers with antibodies	Non-syn SNP/kb PlasmoDB	Non-syn SNPs	
				NF135	NF166
PF3D7_0630600	Conserved hyp	20 (52.6)	2.8	0	1
PF3D7_0906500	Arginase	19 (50.0)	6.5	1	1



PF3D7_1456700	Conserved hyp.	19 (50.0)	2.9	0	0
PF3D7_0719700	40S ribosomal S10	15 (39.5)	0	0	0
PF3D7_1219100	Clathrin heavy chain	12 (31.6)	5.5	0	2
PF3D7_0301700	Hypothetical exp.	10 (26.3)	6.3	0	1
PF3D7_1122700	Conserved hyp.	9 (23.7)	7.8	0	0
PF3D7_1455800	LCCL prot.	6 (15.8)	4.1	0	0

441

442 **S4 Figure: Number of people in whom antibodies against *Plasmodium* protein were detected by**  
443 **protein microarray after CPS immunization**

444 Antibody prevalence does not correlate with the number of Non-synonymous SNPs per kb coding  
445 region of the respective gene (PlasmoDB).

446

447 **S5 Figure: Combined levels of polymorphisms for strains NF135 and NF166 among stage-specific**  
448 **proteins.**

449 sporozoite proteins, selected at various levels of stringency, gametocyte proteins, and the remaining  
450 proteins. Sporozoite enriched or sporozoite specific proteins show relatively high levels of  
451 polymorphisms, while gametocyte proteins are clearly depleted of polymorphisms. Furthermore,  
452 antigenic proteins of either stage are enriched relative to non-antigenic proteins

453

454 **S7 Table: Antibody responses in volunteers after CPS immunization.**

455

456 **Variation in sporozoite proteins among selected *Pf* strains NF135 and NF166**

457 Aside from the criteria that the proteins selected as markers for sporozoite exposure are conserved  
458 in sequenced *P. falciparum* strains in PlasmoDB, we also required them to be conserved in two  
459 strains that have been used in research into heterologous protection after CPS: NF135 and NF166.  
460 We therefore sequenced these strains, as well as NF54 (Methods). The mean coverage of mapped  
461 reads in the coding regions ranged from 28 for NF54 to 44 for NF166 (**S8 Table**). A Phred quality score  
462 cut-off of 30 was applied similar to other studies involving *Plasmodium vivax* and *P. falciparum*  
463 sequencing and variant calling, [58-61]. Manual examination of SNPs still uncovered SNPs with a high  
464 Phred score (> 100) that were polymorphic within a single, in principle haploid genome, possibly  
465 reflecting mapping issues in duplicated regions. Next to a Phred score of 30, we also introduced a  
466 required the presence of a variant nucleotide in at least 75% of the high quality bases. Both criteria  
467 were also applied to the indels. Numbers of called SNPs and indels were roughly similar for NF135  
468 and NF166, reflecting their independent geographic origins. As expected, NF54 showed much lower

469 numbers of called SNPs and indels than the other strains, as it is the parental line from which 3D7 is  
470 derived (S8 Table). As observed in Plasmodb, proteins in NF166 and NF135 showed high levels of  
471 polymorphisms for antibody-binding proteins enriched in sporozoites (Figure S6). Nevertheless, most  
472 of them had few or no variations in the set of proteins that we selected as markers for previous  
473 exposure (Table 1). Lowering the maximum allowed variation in these two strains further, e.g. to zero  
474 polymorphisms for all proteins, did not allow us to select a set of proteins that together elicited  
475 antibodies in all 38 volunteers.

476

477 **S6 Figure: Venn diagram with *Plasmodium falciparum* proteins that elicit antibody responses and**  
478 **have non-synonymous SNPs**

479 SNPs in NF135 (blue) and NF166 (yellow) compared to the reference in NF54/3D7. The grey shaded  
480 area contains proteins without any SNPs, they are hence identical to the 3D7 reference and NF54  
481 strain in both NF135 and NF166. The overlap (light green) shows proteins that have SNPs in both  
482 NF135 and NF166, and 27 proteins (dark green) have the exact same SNPs in both strains.

483

484 **S8 Table: Levels of polymorphism in NF54, NF135 and NF166 relative of the reference**  
485 **strain 3D7.** Information for all annotated *Plasmodium falciparum* proteins and the proteins  
486 considered immunogenic and sequencing depth.

487

488

## 489 Discussion

490

491 The sporozoites stage of the *Plasmodium* life cycle represents the parasite's first interaction with the  
492 human immune system and can be used to effectively vaccinate against infection [11]. The set of genes  
493 and proteins expressed at this stage has been determined in at least 11 studies on *Plasmodium spp.*, of  
494 which nine on *P. falciparum*, creating a conundrum of which dataset to use when studying sporozoite  
495 biology, and when deciding to use multiple datasets, how to integrate the datasets. The individual  
496 sporozoite studies did not focus on determining which expression patterns are specific to the  
497 sporozoite stage, rather they examined the presence of transcripts or proteins. Such data of course  
498 can be integrated by combining them in a relatively straightforward manner as has e.g. been done for  
499 sporozoite RNA expression[62], but that does not address how specific the expression of a gene is to  
500 sporozoites. By combining RNA and protein expression data measured across life stages with sets of  
501 proteins known to be present or absent from sporozoites in a Bayesian manner we have created a  
502 single list of proteins ranked by their overrepresentation in sporozoites relative to other stages.  
503 Despite translational repression, which is expected to reduce the correlation between mRNA and  
504 protein levels, including the mRNA levels led to a better performance at the protein level that only  
505 including proteomics data. Cross validation shows furthermore that the integrated list is better at  
506 separating the gold standard positive and negative data sets from each other than the individual data  
507 sets.

508 In this study, we did not separate datasets derived from oocyst sporozoites (the earliest form of this  
509 stage) and salivary gland sporozoites (the mature form). There may have been subtle difference in  
510 protein expression between sporozoites in the two differing host environments (midgut versus salivary  
511 gland) that were not detected. However, oocyst-sporozoite data represent a minority of the combined  
512 data set (2/26) and are highly correlated with salivary gland-derived sporozoite data (**S1 Figure**). A list  
513 of proteins with its own sporozoite specificity score will be a valuable resource for studying sporozoite  
514 biology and understanding novel protein function. Genetic manipulations of malaria parasites can only  
515 occur during blood stage development, which makes studying proteins that are essential for both  
516 blood and sporozoite stages difficult. Using our list in combination with the recently published  
517 *piggyBac* whole genome mutagenesis study [63], will allow researchers to determine the approach  
518 required for generating a knockout parasite i.e. whether an inducible (for essential blood stage and  
519 low sporozoite specificity score) or straight knockout (for a high sporozoite specificity score) system  
520 should be considered.

521 In our analysis, we found an enrichment in the proteins involved in type II fatty acid synthesis which  
522 consistent with literature [38]. An increased output of lipids from this pathway may feed into the  
523 production of GPI anchors that are made up of different sugar and lipid components. We did observe  
524 a slight enrichment in proteins involved in creating GPI anchors, i.e. the three GPI  
525 mannosyltransferases. Although there is currently no established relationship between type II fatty  
526 acid pathway and synthesis of GPI anchors[22], it may be interesting to pursue it for this stage of the  
527 life-cycle given the importance of CSP – the most abundant protein on sporozoites, that is anchored to  
528 the surface via a GPI anchor.

529  
530 Processes involved in the production, folding and catabolism proteins were under-represented in the  
531 sporozoite specific and enriched list. Similarly, genes involved in metabolism such as those part of  
532 glycolysis and citric acid cycle were also under-represented, suggestive of sporozoites existing in a low  
533 metabolic state. Sporozoites are released in the mosquito's circulation as early as day 12 post blood  
534 meal [3] and make their way to the salivary gland. Nevertheless they are generally less infectious for  
535 liver cells until after a period of maturation within the salivary gland[64]. Although the exact nature of  
536 this maturation is not known, there is evidence that these parasites remain in a latent state until  
537 ejected into the human host[4, 65]. To understand how sporozoites exist in this state of inactivity, we  
538 have identified several interesting targets such as an ortholog of AKT1/2/3, two 14-3-3 proteins, an  
539 ortholog of the raf kinase inhibitor and an ortholog of the quiescent stem cell marker prominin/CD133.

540  
541 Sequence variation among *Plasmodium* strains is pervasive [66], and is possibly responsible for the  
542 limited heterologous protection after CPS vaccination with NF54 and challenge with NF135 and NF166  
543 [67]. Indeed as we have shown here both sporozoite specific proteins and proteins that elicit  
544 antibodies are highly polymorphic, and only a fraction of those are conserved between NF54, NF135  
545 and NF166 (Fig S6, S7). The correlation between immunogenicity and level of sequence conservation  
546 suggest that antigenic drift plays a role in the sequence variation. It is not clear whether antigenic drift  
547 would also be responsible for high variation among sporozoite specific proteins, as we did not observe  
548 a correlation between the Bayesian score of sporozoite specificity and the immunogenicity.  
549 Nevertheless, among the large number of sporozoite proteins that elicit antibodies there are still  
550 proteins that show limited sequence variation and allow selecting of a small set of proteins that are  
551 well conserved and against which antibodies could serve as potential vaccine targets or markers of  
552 previous exposure to sporozoites.

553 In summary, we show a set of previously unidentified sporozoite-specific proteins and assign functions  
554 potentially related to the enduring state of inactivity of the salivary gland sporozoite. We further

555 identify sporozoite-directed humoral immune responses and their potential as functional or diagnostic  
556 responses that can be elucidated in future studies.

557

558

559

560

561

## 562 References

563

- 564 1. Amino, R., et al., *Quantitative imaging of Plasmodium transmission from mosquito to*  
565 *mammal*. Nat Med, 2006. **12**(2): p. 220-4.
- 566 2. Matuschewski, K., et al., *Infectivity-associated changes in the transcriptional repertoire of the*  
567 *malaria parasite sporozoite stage*. J Biol Chem, 2002. **277**(44): p. 41948-53.
- 568 3. Yang, A.S.P., et al., *AMA1 and MAEBL are important for Plasmodium falciparum sporozoite*  
569 *infection of the liver*. Cell Microbiol, 2017. **19**(9).
- 570 4. Porter, R.J., R.L. Laird, and E.M. Dusseau, *Studies on malarial sporozoites. II. Effect of age and*  
571 *dosage of sporozoites on their infectiousness*. Exp Parasitol, 1954. **3**(3): p. 267-74.
- 572 5. Le Roch, K.G., et al., *Discovery of gene function by expression profiling of the malaria parasite*  
573 *life cycle*. Science, 2003. **301**(5639): p. 1503-8.
- 574 6. Florens, L., et al., *A proteomic view of the Plasmodium falciparum life cycle*. Nature, 2002.  
575 **419**(6906): p. 520-6.
- 576 7. Lindner, S.E., et al., *Total and putative surface proteomics of malaria parasite salivary gland*  
577 *sporozoites*. Mol Cell Proteomics, 2013. **12**(5): p. 1127-43.
- 578 8. Bennink, S. and G. Pradel, *The molecular machinery of translational control in malaria*  
579 *parasites*. Mol Microbiol, 2019. **112**(6): p. 1658-1673.
- 580 9. Zimmermann, L., et al., *A Completely Reimplemented MPI Bioinformatics Toolkit with a New*  
581 *HHpred Server at its Core*. J Mol Biol, 2018. **430**(15): p. 2237-2243.
- 582 10. Szklarczyk, R., et al., *Iterative orthology prediction uncovers new mitochondrial proteins and*  
583 *identifies C12orf62 as the human ortholog of COX14, a protein involved in the assembly of*  
584 *cytochrome c oxidase*. Genome Biol, 2012. **13**(2): p. R12.
- 585 11. Obiero, J.M., et al., *Antibody Biomarkers Associated with Sterile Protection Induced by*  
586 *Controlled Human Malaria Infection under Chloroquine Prophylaxis*. mSphere, 2019. **4**(1).
- 587 12. Aurrecochea, C., et al., *PlasmoDB: a functional genomic database for malaria parasites*.  
588 Nucleic Acids Res, 2009. **37**(Database issue): p. D539-43.
- 589 13. Meerstein-Kessel, L., et al., *Probabilistic data integration identifies reliable gametocyte-*  
590 *specific proteins and transcripts in malaria parasites*. Scientific Reports, 2018. **8**: p. 13.
- 591 14. Subramanian, A., et al., *Gene set enrichment analysis: a knowledge-based approach for*  
592 *interpreting genome-wide expression profiles*. Proc Natl Acad Sci U S A, 2005. **102**(43): p.  
593 15545-50.
- 594 15. Lindner, S.E., J.L. Miller, and S.H. Kappe, *Malaria parasite pre-erythrocytic infection:*  
595 *preparation meets opportunity*. Cell Microbiol, 2012. **14**(3): p. 316-24.
- 596 16. Dessens, J.T., et al., *CTRP is essential for mosquito infection by malaria ookinetes*. Embo  
597 Journal, 1999. **18**(22): p. 6221-6227.
- 598 17. Engelmann, S., O. Silvie, and K. Matuschewski, *Disruption of Plasmodium Sporozoite*  
599 *Transmission by Depletion of Sporozoite Invasion-Associated Protein 1*. Eukaryotic Cell, 2009.  
600 **8**(4): p. 640-648.

- 601 18. Boucher, L.E. and J. Bosch, *The apicomplexan glideosome and adhesins - Structures and*  
602 *function*. J Struct Biol, 2015. **190**(2): p. 93-114.
- 603 19. Green, J.L., et al., *Compositional and expression analyses of the glideosome during the*  
604 *Plasmodium life cycle reveal an additional myosin light chain required for maximum motility*.  
605 J Biol Chem, 2017. **292**(43): p. 17857-17875.
- 606 20. Kumar, V., et al., *Inner membrane complex 1l protein of Plasmodium falciparum links*  
607 *membrane lipids with cytoskeletal element 'actin' and its associated motor 'myosin'*. Int J Biol  
608 Macromol, 2019. **126**: p. 673-684.
- 609 21. Gilson, P.R., et al., *Identification and stoichiometry of glycosylphosphatidylinositol-anchored*  
610 *membrane proteins of the human malaria parasite Plasmodium falciparum*. Mol Cell  
611 Proteomics, 2006. **5**(7): p. 1286-99.
- 612 22. Shears, M.J., C.Y. Botte, and G.I. McFadden, *Fatty acid metabolism in the Plasmodium*  
613 *apicoplast: Drugs, doubts and knockouts*. Mol Biochem Parasitol, 2015. **199**(1-2): p. 34-50.
- 614 23. Delemarre, B.J. and H.J. van der Kaay, *[Tropical malaria contracted the natural way in the*  
615 *Netherlands]*. Ned Tijdschr Geneesk, 1979. **123**(46): p. 1981-2.
- 616 24. Teirlinck, A.C., et al., *NF135.C10: a new Plasmodium falciparum clone for controlled human*  
617 *malaria infections*. J Infect Dis, 2013. **207**(4): p. 656-60.
- 618 25. McCall, M.B.B., et al., *Infectivity of Plasmodium falciparum sporozoites determines emerging*  
619 *parasitemia in infected volunteers*. Sci Transl Med, 2017. **9**(395).
- 620 26. Andrews, S., *FastQC: A quality control for high throughput sequence data*. 2010.
- 621 27. Krueger, F., *Trim Galore. A wrapper tool around Cutadapt and FastQC to consistently apply*  
622 *quality and adapter trimming to FastQC files*. 2016.
- 623 28. Coolen, J., *CleanNextSeq\_paired*. 2017.
- 624 29. Langmead, B. and S.L. Salzberg, *Fast gapped-read alignment with Bowtie 2*. Nat Methods,  
625 2012. **9**(4): p. 357-9.
- 626 30. Li, H., et al., *The Sequence Alignment/Map format and SAMtools*. Bioinformatics, 2009.  
627 **25**(16): p. 2078-9.
- 628 31. Quinlan, A.R. and I.M. Hall, *BEDTools: a flexible suite of utilities for comparing genomic*  
629 *features*. Bioinformatics, 2010. **26**(6): p. 841-2.
- 630 32. Lawrence, M., et al., *Software for computing and annotating genomic ranges*. PLoS Comput  
631 Biol, 2013. **9**(8): p. e1003118.
- 632 33. Fidock, D.A., et al., *Cloning and characterization of a novel Plasmodium falciparum sporozoite*  
633 *surface antigen, STARP*. Mol Biochem Parasitol, 1994. **64**(2): p. 219-32.
- 634 34. Gardner, M.J., et al., *Genome sequence of the human malaria parasite Plasmodium*  
635 *falciparum*. Nature, 2002. **419**(6906): p. 498-511.
- 636 35. Arredondo, S.A., et al., *Structure of the Plasmodium 6-cysteine s48/45 domain*. Proc Natl  
637 Acad Sci U S A, 2012. **109**(17): p. 6692-7.
- 638 36. Weigmann, A., et al., *Prominin, a novel microvilli-specific polytopic membrane protein of the*  
639 *apical surface of epithelial cells, is targeted to plasmalemmal protrusions of non-epithelial*  
640 *cells*. Proc Natl Acad Sci U S A, 1997. **94**(23): p. 12425-30.
- 641 37. van Dooren, G.G. and B. Striepen, *The algal past and parasite present of the apicoplast*. Annu  
642 Rev Microbiol, 2013. **67**: p. 271-89.
- 643 38. van Schaijk, B.C., et al., *Type II fatty acid biosynthesis is essential for Plasmodium falciparum*  
644 *sporozoite development in the midgut of Anopheles mosquitoes*. Eukaryot Cell, 2014. **13**(5): p.  
645 550-9.
- 646 39. Bethke, L.L., et al., *Duplication, gene conversion, and genetic diversity in the species-specific*  
647 *acyl-CoA synthetase gene family of Plasmodium falciparum*. Mol Biochem Parasitol, 2006.  
648 **150**(1): p. 10-24.
- 649 40. McConville, M.J. and M.A. Ferguson, *The structure, biosynthesis and function of glycosylated*  
650 *phosphatidylinositols in the parasitic protozoa and higher eukaryotes*. Biochem J, 1993. **294** (  
651 **Pt 2**): p. 305-24.



- 652 41. Englund, P.T., *The structure and biosynthesis of glycosyl phosphatidylinositol protein anchors.*  
653 *Annu Rev Biochem*, 1993. **62**: p. 121-38.
- 654 42. Haldar, K., C.L. Henderson, and G.A. Cross, *Identification of the parasite transferrin receptor*  
655 *of Plasmodium falciparum-infected erythrocytes and its acylation via 1,2-diacyl-sn-glycerol.*  
656 *Proc Natl Acad Sci U S A*, 1986. **83**(22): p. 8565-9.
- 657 43. Ejigiri, I., et al., *Shedding of TRAP by a rhomboid protease from the malaria sporozoite surface*  
658 *is essential for gliding motility and sporozoite infectivity.* *PLoS Pathog*, 2012. **8**(7): p.  
659 e1002725.
- 660 44. Lindner, S.E., et al., *Transcriptomics and proteomics reveal two waves of translational*  
661 *repression during the maturation of malaria parasite sporozoites.* *Nat Commun*, 2019. **10**(1):  
662 p. 4964.
- 663 45. Mair, G.R., et al., *Regulation of sexual development of Plasmodium by translational*  
664 *repression.* *Science*, 2006. **313**(5787): p. 667-9.
- 665 46. Modrzynska, K., et al., *A Knockout Screen of ApiAP2 Genes Reveals Networks of Interacting*  
666 *Transcriptional Regulators Controlling the Plasmodium Life Cycle.* *Cell Host Microbe*, 2017.  
667 **21**(1): p. 11-22.
- 668 47. Petter, M., et al., *Expression of P. falciparum var genes involves exchange of the histone*  
669 *variant H2A.Z at the promoter.* *PLoS Pathog*, 2011. **7**(2): p. e1001292.
- 670 48. Mbengue, A., et al., *A molecular mechanism of artemisinin resistance in Plasmodium*  
671 *falciparum malaria.* *Nature*, 2015. **520**(7549): p. 683-7.
- 672 49. Hoxhaj, G. and B.D. Manning, *The PI3K-AKT network at the interface of oncogenic signalling*  
673 *and cancer metabolism.* *Nat Rev Cancer*, 2020. **20**(2): p. 74-88.
- 674 50. Dastidar, E.G., et al., *Comprehensive histone phosphorylation analysis and identification of*  
675 *Pf14-3-3 protein as a histone H3 phosphorylation reader in malaria parasites.* *PLoS One*,  
676 2013. **8**(1): p. e53179.
- 677 51. Kugelstadt, D., et al., *Raf kinase inhibitor protein affects activity of Plasmodium falciparum*  
678 *calcium-dependent protein kinase 1.* *Mol Biochem Parasitol*, 2007. **151**(1): p. 111-7.
- 679 52. Vandamme, D., et al., *Regulation of the MAPK pathway by raf kinase inhibitory protein.* *Crit*  
680 *Rev Oncog*, 2014. **19**(6): p. 405-15.
- 681 53. Flores-Guzman, F., J. Utikal, and V. Umansky, *Dormant tumor cells interact with memory*  
682 *CD8(+) T cells in RET transgenic mouse melanoma model.* *Cancer Lett*, 2020. **474**: p. 74-81.
- 683 54. Luo, Y., et al., *Single-cell transcriptome analyses reveal signals to activate dormant neural*  
684 *stem cells.* *Cell*, 2015. **161**(5): p. 1175-1186.
- 685 55. Wei, Y., et al., *Activation of PI3K/Akt pathway by CD133-p85 interaction promotes*  
686 *tumorigenic capacity of glioma stem cells.* *Proc Natl Acad Sci U S A*, 2013. **110**(17): p. 6829-  
687 34.
- 688 56. Liu, J., et al., *Novel thioredoxin-like proteins are components of a protein complex coating the*  
689 *cortical microtubules of Toxoplasma gondii.* *Eukaryot Cell*, 2013. **12**(12): p. 1588-99.
- 690 57. Spreng, B., et al., *Microtubule number and length determine cellular shape and function in*  
691 *Plasmodium.* *EMBO J*, 2019. **38**(15): p. e100984.
- 692 58. Diez Benavente, E., et al., *Genomic variation in Plasmodium vivax malaria reveals regions*  
693 *under selective pressure.* *PLoS One*, 2017. **12**(5): p. e0177134.
- 694 59. Shen, H.M., et al., *Genome-wide scans for the identification of Plasmodium vivax genes under*  
695 *positive selection.* *Malar J*, 2017. **16**(1): p. 238.
- 696 60. Campino, S., et al., *Genomic variation in two gametocyte non-producing Plasmodium*  
697 *falciparum clonal lines.* *Malar J*, 2016. **15**: p. 229.
- 698 61. Ocholla, H., et al., *Whole-genome scans provide evidence of adaptive evolution in Malawian*  
699 *Plasmodium falciparum isolates.* *J Infect Dis*, 2014. **210**(12): p. 1991-2000.
- 700 62. Ruiz, J.L. and E. Gomez-Diaz, *The second life of Plasmodium in the mosquito host: gene*  
701 *regulation on the move.* *Brief Funct Genomics*, 2019. **18**(5): p. 313-357.
- 702 63. Zhang, M., et al., *Uncovering the essential genes of the human malaria parasite Plasmodium*  
703 *falciparum by saturation mutagenesis.* *Science*, 2018. **360**(6388).

- 704 64. Ghosh, A.K. and M. Jacobs-Lorena, *Plasmodium sporozoite invasion of the mosquito salivary*  
705 *gland*. *Curr Opin Microbiol*, 2009. **12**(4): p. 394-400.
- 706 65. Beier, J.C., *Malaria parasite development in mosquitoes*. *Annu Rev Entomol*, 1998. **43**: p. 519-  
707 43.
- 708 66. Amambua-Ngwa, A., et al., *Major subpopulations of Plasmodium falciparum in sub-Saharan*  
709 *Africa*. *Science*, 2019. **365**(6455): p. 813-816.
- 710 67. Walk, J., et al., *Modest heterologous protection after Plasmodium falciparum sporozoite*  
711 *immunization: a double-blind randomized controlled clinical trial*. *BMC Med*, 2017. **15**(1): p.  
712 168.
- 713



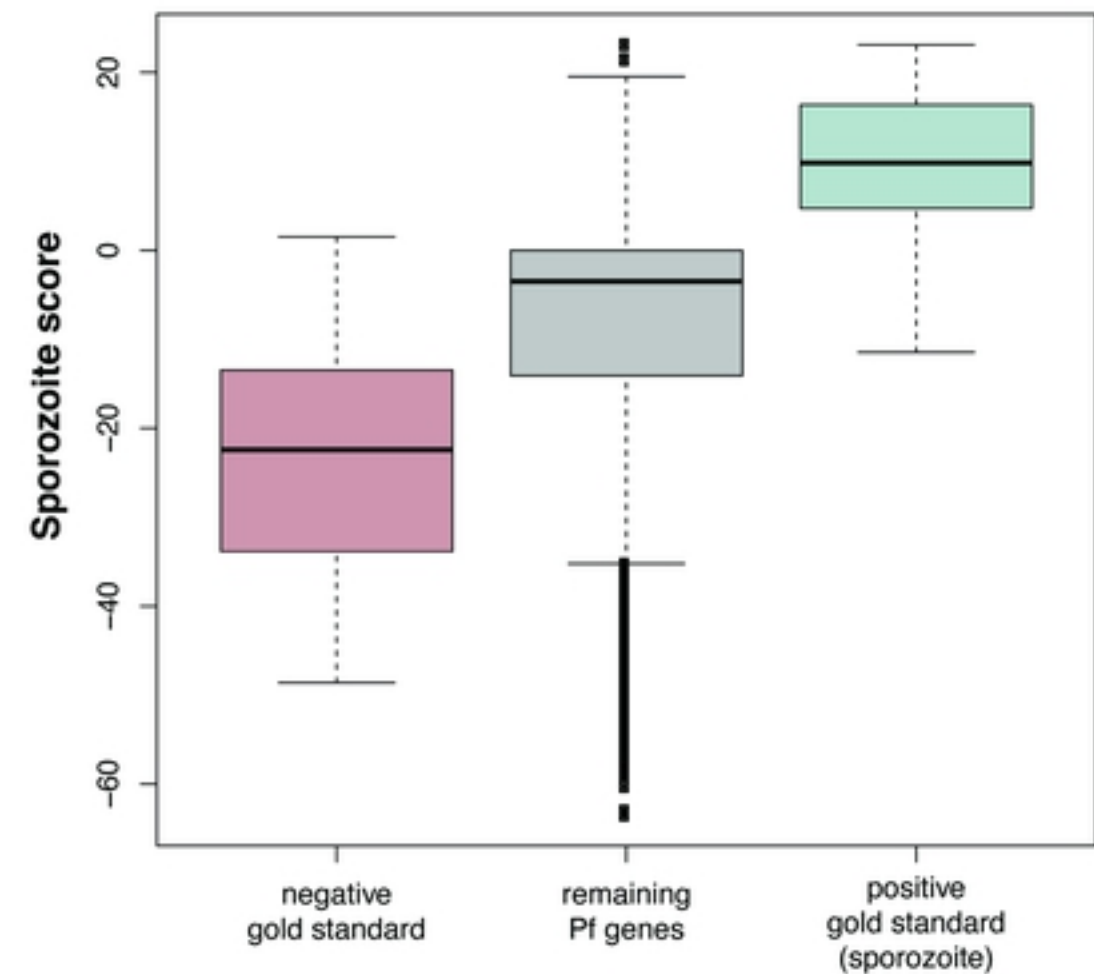
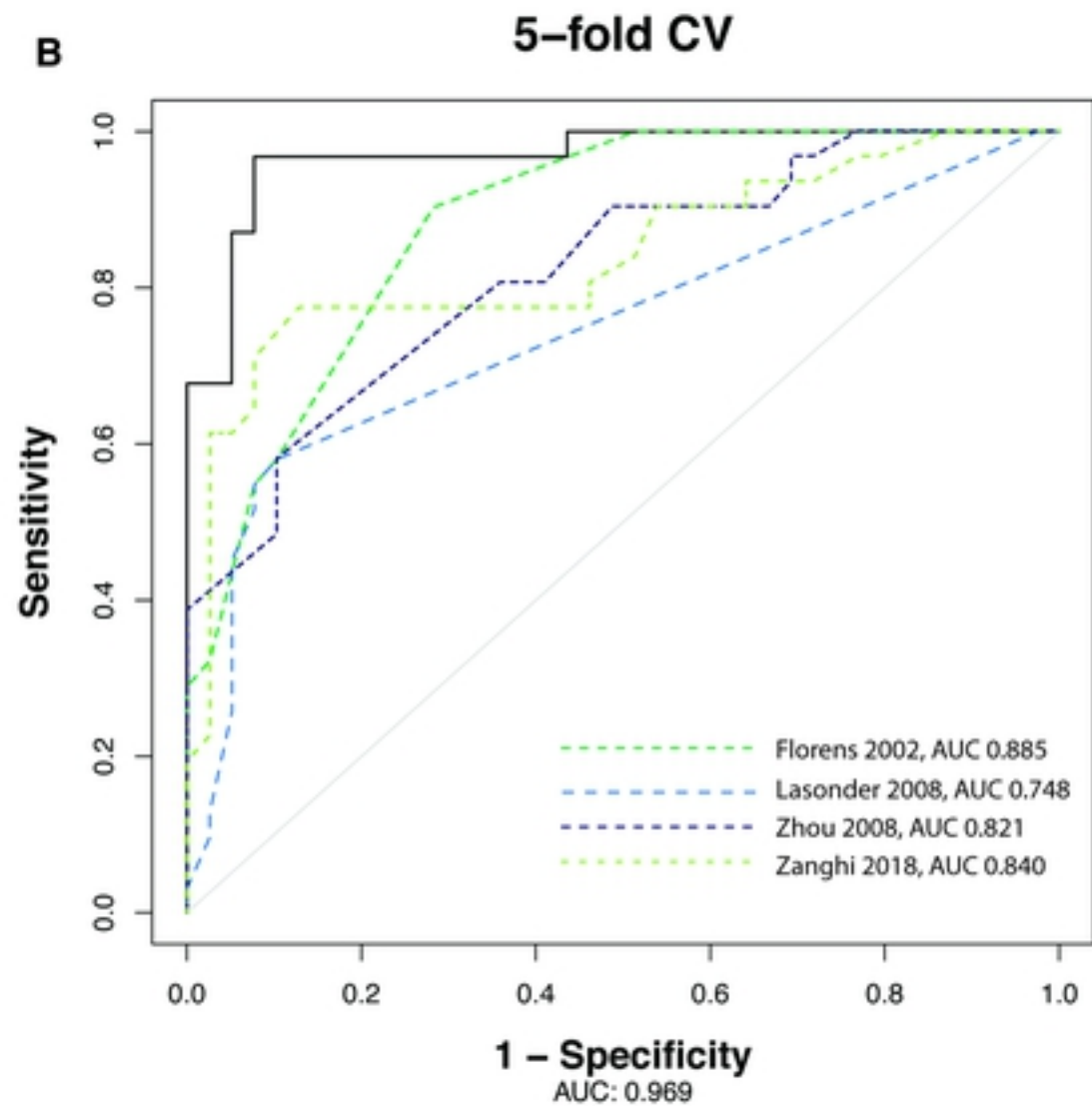
**A****B**

Figure 1

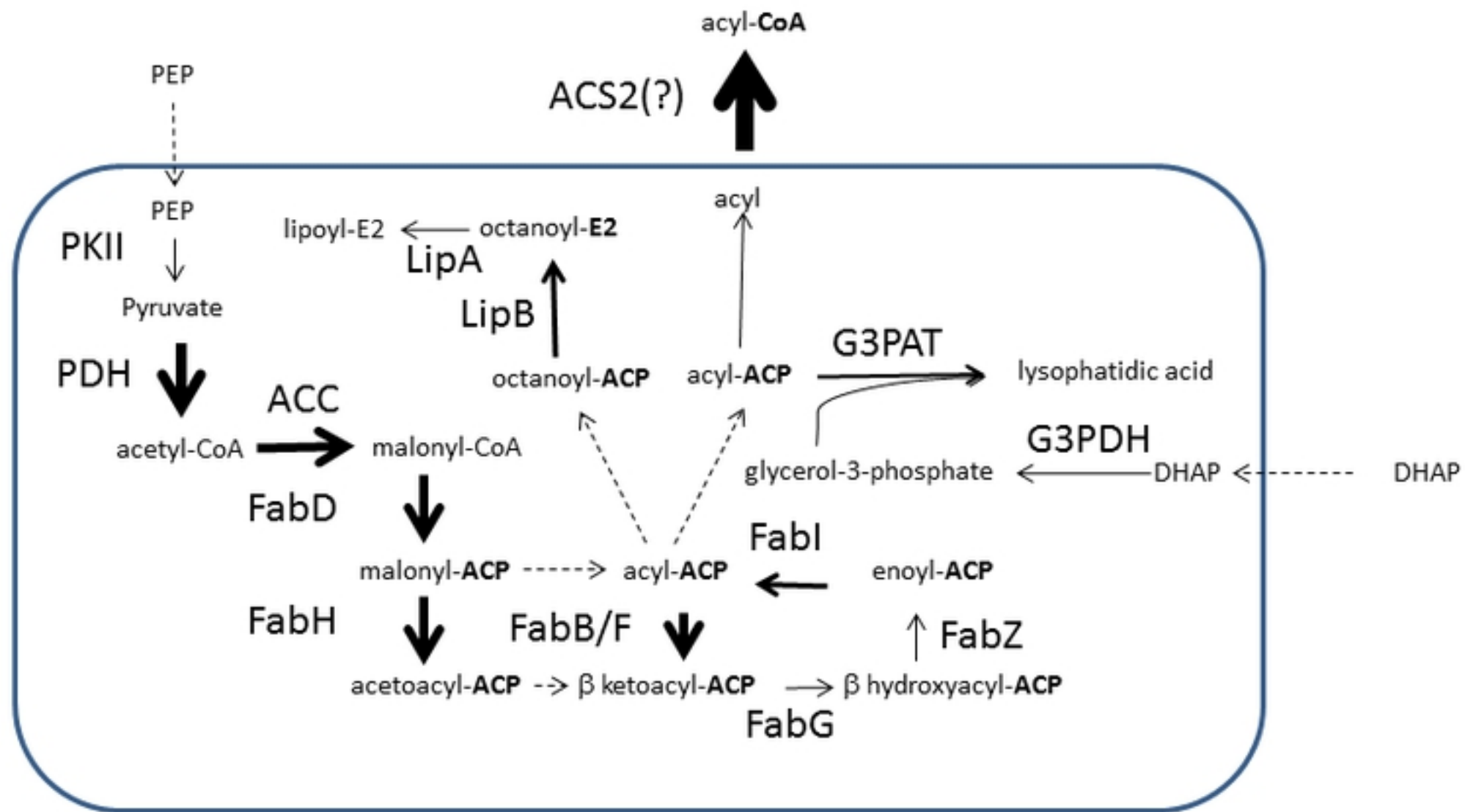


Figure3

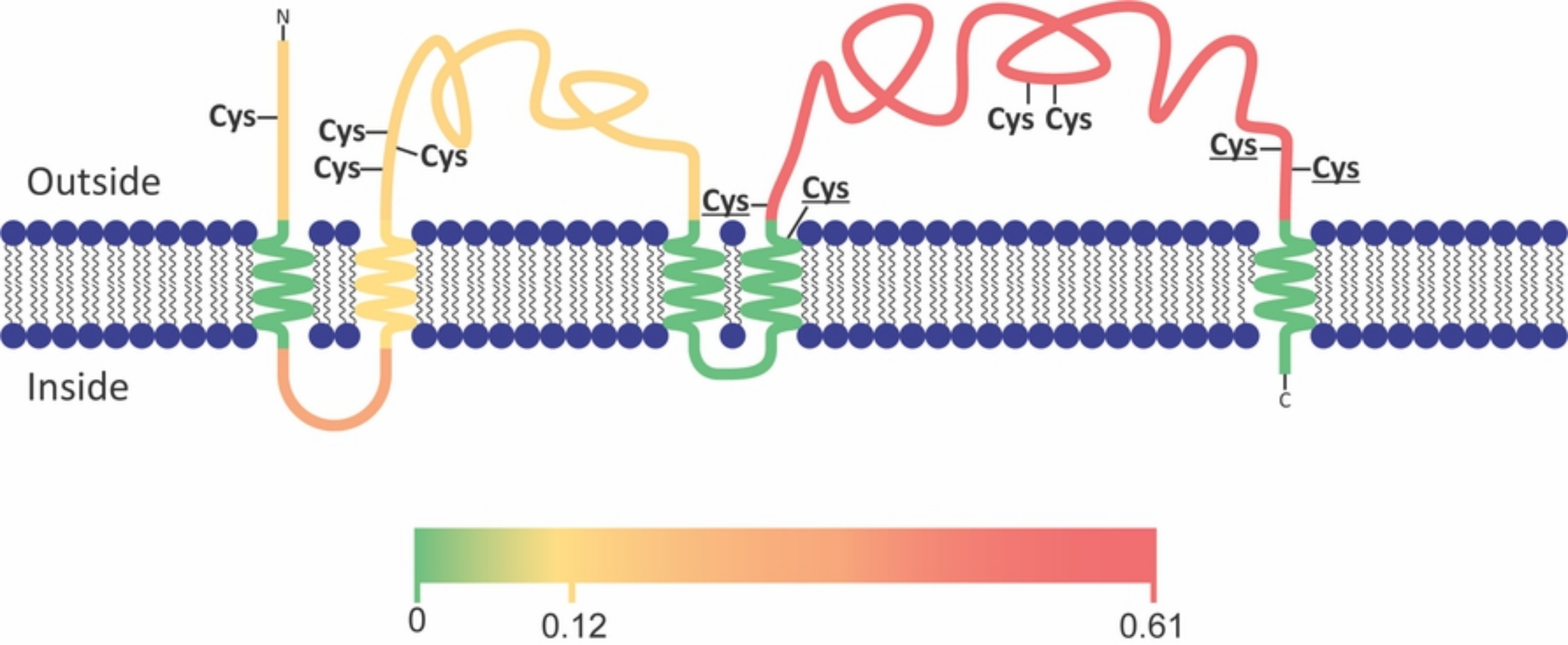


Figure2