1  Phosphorylation in the *Plasmodium falciparum* proteome: A meta-analysis of publicly available data
2  sets

3

4  *Oscar J M Camacho[1], Kerry A Ramsbottom[1], Ananth Prakash[2], Zhi Sun[3], Yasset Perez Riverol[2], Emily*
5  *Bowler-Barnett[2], Maria Martin[2], Jun Fan[2], Eric W Deutsch[3], Juan Antonio Vizcaíno[2] and Andrew R Jones[1]\**

6

7  [1]Institute of Systems, Molecular and Integrative Biology, University of Liverpool, Liverpool, L69 7BE,
8  United Kingdom

9  [2]European Molecular Biology Laboratory, EMBL-European Bioinformatics Institute (EMBL-EBI), Hinxton,
10  Cambridge, CB10 1SD, United Kingdom.

11  [3]Institute for Systems Biology, Seattle, Washington 98109, United States

12  *Andrew.Jones@liverpool.ac.uk

13

14  **Abstract**

15  Malaria is a deadly disease caused by Apicomplexan parasites of the *Plasmodium* genus. Several

16  species of the *Plasmodium* genus are known to be infectious to human, of which *P. falciparum* is the

17  most virulent. Post-translational modifications (PTMs) of proteins coordinate cell signalling and hence,

18  regulate many biological processes in *P. falciparum* homeostasis and host infection, of which the most

19  highly studied is phosphorylation. Phosphosites on proteins can be identified by tandem mass

20  spectrometry (MS) performed on enriched samples (phosphoproteomics), followed by downstream

21  computational analyses. We have performed a large-scale meta-analysis of 11 publicly available

22  phosphoproteomics datasets, to build a comprehensive atlas of phosphosites in the *P. falciparum*

23  proteome, using robust pipelines aimed at strict control of false identifications. We identified a total

24  of 28,495 phosphorylated sites on *P. falciparum* proteins at 5% false localisation rate (FLR) and, of

25  those, 18,100 at 1% FLR. We identified significant sequence motifs, likely indicative of different groups

26  of kinases, responsible for different groups of phosphosites. Conservation analysis identified clusters

27  of phosphoproteins that are highly conserved, and others that are evolving faster within the

28  *Plasmodium* genus, and implicated in different pathways. We were also able to identify over 180,000

29  phosphosites within *Plasmodium* species beyond *falciparum*, based on orthologue mapping.

30   We also explored the structural context of phosphosites, identifying a strong enrichment for

31   phosphosites on fast evolving (low conservation) intrinsically disordered regions (IDRs) of proteins. In

32   other species, IDRs have been shown to have an important role in modulating protein-protein

33   interactions, particularly in signalling, and thus warranting further study for their roles in host-

34   pathogen interactions. All data has made available via UniProtKB, PRIDE and PeptideAtlas, with

35   visualisation interfaces for exploring phosphosites in the context of other data on *Plasmodium*

36   proteins.

37

38   **Author Summary**

39   *Plasmodium* parasites continue to pose a significant global health threat, with a high proportion of the

40   world at risk of malaria. It is imperative to gain new insights into cell signalling and regulation of

41   biological processes in these parasites to develop effective treatments. This study focused on post-

42   translational modifications (PTMs) of proteins, specifically phosphorylation. We conducted a meta-

43   analysis of 11 publicly available phosphoproteomics datasets, identifying over 28,000 phosphorylated

44   sites on *P. falciparum* proteins, using very rigorous statistics to avoid reporting false positives, and

45   mapping to over 180,000 phosphorylation sites on other species of *Plasmodium*.

46   The analysis revealed distinct sequence motifs associated with different groups of phosphosites (and

47   likely indicative of different upstream kinases), and differences in the downstream pathways

48   regulated. Conservation analysis highlighted clusters of phosphoproteins evolving at different rates

49   within the *Plasmodium* genus. Notably, phosphorylation was enriched in regions of proteins lacking

50   distinct structural elements, known as intrinsically disordered regions (IDRs), which are poorly

51   conserved across the genus – we speculate that they are important for modulating protein

52   interactions. The findings provide valuable insights into the molecular mechanisms of *P. falciparum*,

53   with potential implications for understanding host-pathogen interactions. The comprehensive dataset

54   generated is now publicly accessible, serving as a valuable resource for the scientific community

55   through UniProtKB, PRIDE, and PeptideAtlas.

56

**Introduction**

57

58　Malaria remains a major global health burden with 247 million cases worldwide in 2021. In the same

59　year, the World Health Organisation (WHO) has been estimated that 619,000 people died from the

60　disease. Most malaria cases (95%) and deaths (96%) occurred in Sub-Saharan Africa. Malaria is caused

61　by apicomplexan parasites of the *Plasmodium* genus. Of the approximately 156 named *Plasmodium*

62　species, only five have been found to infect humans: *P. falciparum*, *P. vivax*, *P. ovale*, *P. malariae*, and

63　*P. knowlesi*. *Plasmodium* is transmitted from one human to another by female *Anopheles* mosquitos

64　with the exception of *P. knowlesi*, which is believed to be zoonotic i.e. transmission happens from

65　macaques to humans in southeast Asia, where macaques have been previously infected by *Anopheles*

66　mosquitos [1]. Most severe cases and deaths from malaria are caused by *P. falciparum* infections,

67　endemic to Sub-Saharan Africa.

68　The *P. falciparum* life cycle requires two hosts, an *Anopheles* mosquito (around 40 species of

69　*Anopheles* can transmit *P. falciparum* [2]) and a human host. Extracellular sporozoites are transmitted

70　to the human dermal tissue during a blood meal. Once they reach the liver they replicate and develop

71　into merozoites and get released into the peripheral blood. There, merozoites break into erythrocytes

72　and replicate, causing most malaria symptoms. A small proportion of parasites will develop into

73　gametocytes with sexual attributes, which, closing the cycle, are transmitted back to mosquitoes

74　where sporozoites are formed from gametocytes [3].

75　Post-translational modifications (PTMs) of proteins are critically important as they can act as

76　molecular switches. PTMs and specifically phosphorylation has been shown to dynamically change,

77　for example, between extra and intraerythrocytic life-cycle stages suggesting protein sets and

78　pathways with roles in cell invasion [4]. Across stages in the *Plasmodium* intraerythrocytic asexual

79　cycle, a study has reported changes in abundance, defined as peak changes of at least 1.5 fold between

80　stages, for 34% of identified proteins and 75% of phosphorylation sites [5]. Interruption and

81    obstruction of interactions among these proteins could constitute efficient treatments against

82    malaria.

83    All stages in the *Plasmodium* life cycle have the potential to generate targets for vaccines. For example,

84    transmission-blocking vaccines preventing mosquito infection or interfering in *Plasmodium* sexual

85    stages; pre-erythrocytic stage by disabling the ability of merozoites to reach the human liver or

86    replicate there; or targeting interactions at the blood stage by suppressing their ability to enter

87    erythrocytes or replicate [6]. The most advanced vaccine (RTS,S), targets the *P. falciparum*

88    circumsporozoite surface protein (PfCSP) [7-9]. Artemisinin-based combination treatments (ACTs)

89    have proved effective for treating *P. falciparum* malaria [10, 11] but increasing *Plasmodium* drug

90    resistance has been observed, highlighting the importance of development of new drugs [12].

91    There are several useful online resources to support research in *Plasmodium*. PlasmoDB [13, 14]

92    stands out, it is part of the Eukaryotic Pathogen, Vector and Host Informatics Resources [15]

93    (VEuPathDB), which has been running since 2004, collating genome, functional genomic and

94    phenotypic data sets for multiple *Plasmodium* species. PlasmoDB hosts 20 *Plasmodium falciparum*

95    annotated genomes, including the canonical reference P. *falciparum* clone 3D7 [16]. A search of the

96    *P. falciparum* proteome in the UniProt Knowledgebase (UniProtKB) [17], the world's most popular

97    protein knowledge-base, returns 18 proteomes linked to different isolates of countries of origin, which

98    are mostly well synchronised with PlasmoDB.

99    Tandem mass spectrometry (MS/MS) is most often used in large scale phosphosite identification and

100   localisation studies [18]. Protein samples are purified and enzymatically digested, typically using

101   trypsin. Samples are then enriched for phosphorylated peptides using reagents such as $TiO_2$, or other

102   metal ions that promote phosphate binding. Then liquid chromatography is used to separate peptides

103   that are subsequently fragmented and analysed by MS/MS. Results from MS analyses can then be

104   compared against protein sequence databases, with and without the mass shift for phosphorylation,

105   via one of the many available search algorithms [19, 20]. Algorithms provide identification of peptides

106    and localisation of PTM sites in those peptides with scores aiming to reflect the level of confidence

107    that those identifications are correct. Score thresholding is used to select a subset of what are

108    expected to be the most confident findings. However, score thresholding does not provide

109    information about the global false discovery rate (FDR) of peptides, or the global false localisation rate

110    (FLR) of the phosphosites within those peptides. Absence of objective calculation of FLR in phophosite

111    localisation studies hinders comparisons among studies, as it is not possible to establish a common

112    quality threshold among results from different studies. To overcome this problem, we have recently

113    published an approach that allows estimation of global site-level FLR, by including a decoy amino acid,

114    specifically Alanine, for phosphorylation searches (which cannot be modified) as a search parameter

115    to compete against targets sites (S, T or Y), i.e. the pASTY method [21]. An important benefit of pASTY

116    searches is that it allows combination of results from multiple studies as FLR can provide objective

117    comparable thresholds [22].

118    For *P. falciparum,* PlasmoDB provides information on 16,118 phosphorylation sites, although

119    phosphorylation sites have been loaded from multiple publications over the last 10 years. In our

120    previous work examining databases containing human phosphosites, we estimated that there is a high

121    proportion of false positive sites recorded, due to historically inadequate statistics associated with

122    detection of sites by MS [23] and, prior to [22], lack of methods for calculating adequate statistics for

123    controlling the FLR across studies.

124    In this work, we aim to provide a high-quality mapping of *P. falciparum* phosphosites via a large-scale

125    re-analysis of public phospho-enriched studies, underpinned by robust analysis pipelines enabling

126    meta-analyses with FLR control within and across studies' results. This analysis is part of the

127    "PTMeXchange" initiative, which is re-analysing phospho-enriched data sets and depositing results

128    into the proteomics resources PRIDE [24], PeptideAtlas [25] and UniProtKB. Results from downstream

129    analysis in the most confident phosphosites are also reported in this manuscript, including analysis

130    motifs centred on phosphosites and pathway enrichment analysis for these motifs. We examine

131 phosphorylation site conservation between our reference isolate 3D7 (*P. falciparum*) and other

132 species of the *Plasmodium* genus, as well as investigating the structure and disordered regions of

133 phosphoproteins.

**Results**

**Phosphosite identification**

136 The counts of peptide-spectrum matches (PSMs), and PSM-sites passing the 1% FDR threshold are

137 displayed in Table 1 as well as the number of overall PSM-sites at 1% and 5% FLR. Next, data were

138 collapsed to peptidoform-site level *i.e.* removing redundancy caused by the common occurrence of

139 multiple PSMs identifying the same peptidoform. A peptidoform is defined as a unique sequence of

140 amino acids with specific modifications. For example, two identical peptide sequences but with

141 modifications at different positions in the sequence are different peptidoforms. Table 1 also displays

142 the number of peptidoform-sites and protein-sites (accepting the mapping from a peptide sequence

143 to all proteins it can be found in, assuming tryptic cleavage) at 1% and 5% FLR, with these counts

144 separated for *P. falciparum* and human matches. Note that FLR threshold counts only considered

145 PSMs that passed 1% FDR. The data is unequally distributed between the 11 studies, with four studies

146 (PXD012143, PXD015833, PXD020381, PXD026474) contributing significantly more to the overall

147 number of sites at every stage of the analysis. In fact, nearly 70% of all *P. falciparum* protein-sites at

148 5% FLR come from these four studies.

149 *Table 1. From left to right: PSM count at 1% FDR, PSM-site count at 1% FDR, overall PSM-site count at 1% and 5% FLR, P.*

150 *falciparum peptidoform-site count at 1% and 5% FLR, human peptidoform-site count at 1% and 5% FLR, P. falciparum protein-*

151 *site count at 1% and 5% FLR, human protein-site count at 1% and 5% FLR.*

| | | | Overall | | Plasmodium | | Human | | Plasmodium | | Human | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1% FDR PSM | 1% FDR PSM-site | 1% FLR PSM-Site | 5% FLR PSM-Site | 1% FLR Peptidoform-site | 5% FLR Peptidoform-site | 1% FLR Peptidoform-site | 5% FLR Peptidoform-site | 1% FLR Protein-site | 5% FLR Protein-site | 1% FLR Protein-site | 5% FLR Protein-site |
| PXD000070 | 42054 | 46064 | 6206 | 12375 | 564 | 1265 | 44 | 103 | 398 | 886 | 27 | 69 |
| PXD001684 | 14152 | 14179 | 5648 | 12015 | 901 | 1291 | 48 | 56 | 785 | 1110 | 45 | 63 |
| PXD002266 | 38618 | 45004 | 10342 | 18609 | 1609 | 3298 | 68 | 150 | 1147 | 2421 | 47 | 118 |
| PXD005207 | 21967 | 24045 | 4584 | 8105 | 2367 | 3386 | 188 | 296 | 1364 | 1884 | 105 | 174 |
| PXD009157 | 22449 | 24653 | 4134 | 7049 | 2179 | 3353 | 83 | 151 | 1291 | 2023 | 31 | 69 |

| PXD009465 | 24435 | 37096 | 7812 | 10694 | 3905 | 5382 | 101 | 159 | 2513 | 3534 | 70 | 122 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| PXD012143 | 270306 | 340996 | 93438 | 144256 | 26893 | 39487 | 633 | 1139 | 10597 | 15797 | 350 | 662 |
| PXD015093 | 28693 | 30469 | 5597 | 12187 | 1156 | 2760 | 123 | 274 | 583 | 1242 | 64 | 149 |
| PXD015833 | 348842 | 392870 | 78345 | 134576 | 22288 | 43134 | 1291 | 3005 | 6439 | 11897 | 442 | 1060 |
| PXD020381 | 228888 | 278203 | 67457 | 105938 | 17604 | 31622 | 361 | 832 | 7331 | 13050 | 173 | 483 |
| PXD026474 | 154644 | 168133 | 57786 | 100714 | 8772 | 14611 | 3334 | 5640 | 4657 | 7426 | 1549 | 2657 |

152

153 Figure 1A shows the Gold-Silver-Bronze (GSB, see Methods) quality categorisation for *P. falciparum*

154 protein-sites matching to *P. falciparum* (left panel) and human (right panel) proteins. If a sequence,

155 and therefore the sites within the sequence, match to more than one protein sequence, these sites

156 were counted for all mapping proteins in the proteome. In total, there were 28,495 protein-sites in *P.*

157 *falciparum* classified as GSB (5% FLR) of which 9,587 were Gold (seen in 2 or more studies at <1% FLR),

158 8,513 Silver (1 study at <1% FLR), and 10,395 Bronze (>=1% FLR and <5% FLR). While for Human (Figure

159 1A right), there were 4,239 protein-sites identified, with 446 Gold, 1,788 Silver and, 2,005 classified as

160 Bronze. 25,758 and 3,529 sites were unique (matching to one protein only), and 351 and 284 non-

161 unique sites for *P. falciparum* and Human respectively, considering only sites at 5% FLR (Figure 1B).

162 The number of GSB unique protein-sites were 26,028 in *P. falciparum* and 3,732 in Human (Figure 1C),

163 with 81 protein-sites matching to both species (this calculation only considers a protein match per

164 sequence and site). Only 1.35% of all protein-sites in our *Plasmodium* analysis mapped to more than

165 1 protein with less than 0.05% mapping to more than two proteins. Interestingly, there was a

166 peptidoform mapping the same site evidence to 36 proteins from the PfEMP1 family, due to very high

167 sequence similarity amongst this gene family (Figure 1D). All Gold, Silver and Bronze phosphosites can

168 be found in Supp File 1.

169

170 Figure 1. A: The count of protein-sites classified as Gold-Silver-Bronze for *P. falciparum* (left) and Human (right), sites

171 coloured by phosphorylated amino acid. This includes all potential locations for the identified sites when peptides match to

172 more than one location or protein and decoy matches to Alanine. B: Number of unique or not unique protein-sites

173 identified within each species among sites at 5% FLR, where not unique are those peptidoforms that map to more than one

174    protein. C: Number of protein-sites, where sites are mapped to a single protein, for each species and common to both. D:

175    Proportion of sites matching to more than one protein.

176    **Motif and pathway enrichment analysis**

177    We investigated if there were overrepresented sequence motifs centred on 5% FLR phosphosites.

178    Motif analysis was performed for *P. falciparum* by comparing 15mer peptides centred on S, T and Y

179    phosphosites against a background of 15mer peptides centred on all STY sites, phosphorylated or not.

180    The analysis returned 107 statistically significant motifs, of which 65, 30 and 12 were centred on S, T

181    and Y, respectively (Supp Figure 1 and Supp File 1).

182    The most common significant motifs (group 1 on Figure 2) were those with S and T combining with E

183    and D in different positions such: [ST]D, D[ST], [ST].[DE], [ST]…D, and [ST]N[DE], N[ST].[DE] or N.[ST]

184    [DE]. There were also other common motifs to S and T not containing D or E like N[ST], K..[ST], R..[ST].

185    Among other potentially interesting features of these motifs, only K was found in motifs at positions -

186    7, -6, +6 and +7, with significant motifs for K…..DS, K….DS and SD….K, SD…..K, SN…..K. Motifs with K

187    relatively distant to the target site may be due to a preference for particular kinases or an artefact

188    related to tryptic cleavage enriching for lysines to be present somewhere in many detected peptides.

189    There were only six 3-amino acid motifs with a P in the +1 position, which can be very common for

190    other species. Many of these motifs agreed with those found by Peace *et al.* [5], which can be expected

191    as their MS data was included in this analysis. Treeck *et al.* [26] also found many of these motifs in *P.*

192    *falciparum* (S[DE].E,SD.[ED], K..S.D, [KR]..S, S[DE], SN, DS).

193    Motifs were independently analysed to investigate functional enrichment of the proteins in which

194    specific motifs were found. Summing across all motif results, 39 different GO terms were found to be

195    statistically significant. Table 3 contains significant GO terms for motifs at least 4-fold enriched versus

196    the background, and the full analysis for all motifs can be found in Supp File 2. It is worth noting that

197    the RSF.D motif gave highly significant matches to several GO terms. However, this result is an unusual

198    artefact of the protocol. *P. falciparum* has a 65 gene family, in which all members are called "PfEMP1",

199    with highly related protein sequences. Phosphosites are mapped to all positions where a peptidoform

8

200     can be found, typically resulting in the vast majority of sites mapping to a single protein (counts for

201     the few exceptions are shown in Figure 1D). Phosphosites identified in PfEMP1 mapped to 36 different

202     proteins, which gives an apparently extremely significant signal under motif-GO analysis, since the

203     PfEMP1 proteins all contain the same motif and are also all mapped to the same GO terms.

204     Table 3. Statistically significant GO terms for the subset of motifs surrounding phosphorylated sites. Analyses were carried

205     out independently for each motif and includes only motifs with at least 4 fold change over background

| ID | Description | GeneRatio | BgRatio | pvalue | p.adjust | qvalue | Count | motif |
|---|---|---|---|---|---|---|---|---|
| GO:0034399 | nuclear periphery | 4/52 | 35/3454 | 0.0016 | 0.0374 | 0.0367 | 4 | .......SD.E..K. |
| GO:0005515 | protein binding | 13/52 | 354/3454 | 0.0017 | 0.0374 | 0.0367 | 13 | .......SD.E..K. |
| GO:0045178 | basal part of cell | 4/144 | 11/3454 | 0.0007 | 0.0314 | 0.0302 | 4 | ....R.NS....... |
| GO:0048870 | cell motility | 4/144 | 12/3454 | 0.0011 | 0.0314 | 0.0302 | 4 | ....R.NS....... |
| GO:0009410 | response to xenobiotic stimulus | 14/144 | 145/3454 | 0.0024 | 0.0470 | 0.0451 | 14 | ....R.NS....... |
| GO:0000792 | heterochromatin | 5/233 | 11/3454 | 0.0004 | 0.0209 | 0.0197 | 5 | .......SD.E.... |
| GO:0005515 | protein binding | 40/233 | 354/3454 | 0.0005 | 0.0209 | 0.0197 | 40 | .......SD.E.... |
| GO:0034399 | nuclear periphery | 8/233 | 35/3454 | 0.0018 | 0.0472 | 0.0446 | 8 | .......SD.E.... |
| GO:0050839 | cell adhesion molecule binding | 30/36 | 54/3454 | 5.59E-53 | 6.71E-52 | 2.35E-52 | 30 | ......RSF.D.... |
| GO:0098609 | cell-cell adhesion | 30/36 | 57/3454 | 5.57E-52 | 3.34E-51 | 1.17E-51 | 30 | ......RSF.D.... |
| GO:0020030 | infected host cell surface knob | 30/36 | 61/3454 | 9.17E-51 | 3.67E-50 | 1.29E-50 | 30 | ......RSF.D.... |
| GO:0020002 | host cell plasma membrane | 30/36 | 214/3454 | 1.10E-31 | 3.30E-31 | 1.16E-31 | 30 | ......RSF.D.... |
| GO:0020013 | modulation by symbiont of host erythrocyte aggregation | 29/36 | 191/3454 | 2.49E-31 | 5.98E-31 | 2.10E-31 | 29 | ......RSF.D.... |
| GO:0020035 | adhesion of symbiont to microvasculature | 28/36 | 172/3454 | 7.88E-31 | 1.58E-30 | 5.53E-31 | 28 | ......RSF.D.... |
| GO:0020033 | antigenic variation | 29/36 | 202/3454 | 1.40E-30 | 2.41E-30 | 8.44E-31 | 29 | ......RSF.D.... |
| GO:0046789 | host cell surface receptor binding | 7/36 | 37/3454 | 5.95E-08 | 8.93E-08 | 3.13E-08 | 7 | ......RSF.D.... |
| GO:0043565 | sequence-specific DNA binding | 4/77 | 17/3454 | 0.0004 | 0.0178 | 0.0160 | 4 | ....N..SP...... |
| GO:0009410 | response to xenobiotic stimulus | 10/75 | 145/3454 | 0.0009 | 0.0484 | 0.0416 | 10 | ....K..T.D..... |
| GO:1903561 | extracellular vesicle | 7/57 | 91/3454 | 0.0006 | 0.0221 | 0.0206 | 7 | ....E..T.E..... |
| GO:0042393 | histone binding | 2/19 | 12/3454 | 0.0018 | 0.0274 | 0.0212 | 2 | .......TP...E.. |
| GO:1903561 | extracellular vesicle | 20/209 | 91/3454 | 2.51E-07 | 1.88E-05 | 1.72E-05 | 20 | ....R..T....... |
| GO:0003723 | RNA binding | 27/209 | 199/3454 | 4.23E-05 | 0.0014 | 0.0013 | 27 | ....R..T....... |
| GO:0020020 | food vacuole | 17/209 | 98/3454 | 5.83E-05 | 0.0014 | 0.0013 | 17 | ....R..T....... |
| GO:0045178 | basal part of cell | 5/209 | 11/3454 | 0.0002 | 0.0049 | 0.0045 | 5 | ....R..T....... |
| GO:0048870 | cell motility | 5/209 | 12/3454 | 0.0004 | 0.0064 | 0.0059 | 5 | ....R..T....... |
| GO:0005515 | protein binding | 35/209 | 354/3454 | 0.0018 | 0.0236 | 0.0216 | 35 | ....R..T....... |
| GO:0042393 | histone binding | 3/30 | 12/3454 | 0.0001 | 0.0033 | 0.0028 | 3 | .....D.TE...... |
| GO:0032991 | protein-containing complex | 2/23 | 14/3454 | 0.0036 | 0.0446 | 0.0405 | 2 | .......Y.SD.... |
| GO:0006511 | ubiquitin-dependent protein catabolic process | 2/23 | 15/3454 | 0.0042 | 0.0446 | 0.0405 | 2 | .......Y.SD.... |
| GO:0003724 | RNA helicase activity | 2/23 | 18/3454 | 0.0061 | 0.0446 | 0.0405 | 2 | .......Y.SD.... |
| GO:0003723 | RNA binding | 6/20 | 199/3454 | 0.0006 | 0.0153 | 0.0147 | 6 | ......SYE...... |

206

207     Under the hypothesis that similar motifs may indicate proteins being involved in similar processes, we

208     identified six groups based on the amino acids forming those motifs (Figure 2). Our analysis showed

209     that besides general agreement on high-level functional terms such as "protein binding", "cytoplasm"

210     or "RNA binding" there were also group specific GO terms pointing to some functional specificity

211     among motifs' groups. There was almost complete agreement in significant terms between groups 1

212     to 3 which could suggest there is overlap in the signalling cascades with regards to downstream

213     effects, or the method may not be sensitive enough to draw more specific terms allowing

214     differentiation. Groups 4 to 6 returned more diverse GO terms, with group 6 likely conditioned by the

215     PfEMP1 protein family as previously explained.

216     As an alternative and more objective approach for grouping motifs, we took a subset of motifs with a

217     4-fold change over background, *i.e.* the most strongly over-represented motifs, and performed an

218     enrichment analysis for all *P. falciparum* proteins where those phosphorylation motifs can be found.

219     Figure 3A shows a heatmap of the p-values for the GO terms associated with proteins containing these

220     phosphorylation motifs. These motifs did not seem to form strong clusters according to these terms,

221     although some small groups (similar pairs clustering together) could be observed. In a few cases, there

222     was some overlap in protein sets carrying motifs, for example SD.E and SD.E..K – the latter being

223     matched to a subset of the proteins matched by the former. However, the motifs K..T.E and K..T.D are

224     mutually exclusive, yet proteins carrying these phosphosite motifs can be seen to be acting in many

225     of the same signalling pathways. Other examples of "pairing" on the dendrogram are motifs N.SP and

226     N..SP, and E..T.E and E...TE. It is possible that these pairs of phospho-motifs are recognised by the

227     same kinase, or that there are closely related kinases within the same family that are involved within

228     the same types of downstream pathways. Another "pair" of phospho-motifs is K..TP and TDD – this

229     latter example is surprising, since it would typically be assumed that S/TP and S/TD phosphosites are

230     regulated by different kinase families. It is possible that there is crosstalk between two kinases,

231     although the evidence here is not sufficient to form any strong conclusions.

232 Therefore, in general, the heatmap suggests some differences in functionality between proteins

233 containing different phosphorylation motifs. However, when analysed as a group, 10 GO terms

234 summarised the functional processes for the subset of genes where these phosphorylated motifs were

235 found (Figure 3B).

236

237 Figure 2. Statistically significant GO terms for motifs surrounding phosphorylated sites, grouped by similar motifs, formed by

238 similar groups of amino acids.

239

240

241 Figure 3. A: Heatmap of the adjusted p-values resulting from an enrichment analysis of the proteins where phosphorylation

242 motifs (on the x-axis) were found (only motifs with fold change over background above 4 are considered). B: Significant GO

243 Terms (p-adjusted <=0.05) for the genes in which motifs with fold change over background above 4 could be found.

244 The *P. falciparum* proteome has around 90 protein kinases – 89 are annotated in PlasmoDB with the

245 InterPro protein kinase domain term. An analysis by Adderley *et al.* [27] identified 98 protein kinases

246 in isolate 3D7, by including keyword searching in addition to InterPro domain searching. From their

247 list of kinases, there several sequences annotated as pseudokinases or having kinase-like domains

248 (PF3D7_0424700, PF3D7_0708300, PF3D7_0724000, PF3D7_0823000, PF3D7_1106800,

249 PF3D7_1321100, PF3D7_1428500) and two pseudogenes (PF3D7_0731400 and PF3D7_1476400).

250 Adderley *et al.* [27] classified these kinases into families, using an HMM (Hidden Markov Models)-

251 based technique. Supp Table 2 shows the counts of 3D7 kinases classified into families, with additional

252 notes on potential kinase motifs for these families based upon information on kinase motifs found in

253 humans from [28], with the caveat that we cannot be sure that even if orthologous kinases exist

254 between humans and *P. falciparum* that they recognise the same motifs. Accurate prediction of

255 kinase-substrate relationships is not straightforward without high-quality experimental data, which is

256 lacking in *P. falciparum.* FIKK kinases are unique to Apicomplexa [29], and seem to have an important

257    role in host invasion (such as phosphorylation of erythrocyte proteins), but otherwise little is known

258    about the motifs for their targets. The publication associated with dataset PXD015833 specifically

259    investigated substrate specificity and concluded that some of the FIKK family have a preference for a

260    basic motif near to pS/pT site – with arginine enriched in the minus position (group 6 in our analysis)

261    [30]. We could speculate that the motifs groups identified in Figure 2 are driven by different kinase

262    groups in Supp. Table 2, but without new experimental data, robust conclusions are not possible. This

263    is an area requiring further work.

264

265

266 **Conservation analysis**

267 Based on unique phosphorylated protein-sites at 5% FLR in *P. falciparum* 3D7, we investigated the

268 existence of their orthologs in other species of *Plasmodium*. Conservation for each site, defined as the

269 proportion of species containing the same amino acid at that position in the multiple sequence

270 alignment, can be found as supplementary files (Supp File 3.). We generated a heatmap (Figure 4)

271 representing the average agreement among sites within each protein. If the site was conserved with

272 respect to the reference 3D7 it was labelled as 1 and 0 otherwise. For each species, a non-conserved

273 site could be the result of having a different amino acid with respect to the reference phosphorylation

274 position, there could be gap in the protein or the orthologue was not found for that species. Then,

275 within each protein and species, the average (proportion of sites conserved) was calculated. In this

276 way, proteins without any sites conserved or when the protein is not found in the species have score=0

277 and when all phosphosites are conserved, the score = 1. Changes in a site between Ser, Thr or Try was

278 considered not conserved, although scores were also calculated allowing for S <-> T substitutions as

279 "conserved" (not disruptive to S/T phosphorylation), as shown in Supp File 3.

280 Mapping across species returned 26,316 sites belonging to 2,890 proteins, where only 1 mapping

281 occurrence was used per phosphorylated protein-site, i.e. when phosphorylated peptides could match

282 to more than one protein (or very rarely multiple positions within one protein) only one match was

283 used per peptide. Of those 2,890 proteins, only 108 proteins contained all phosphosites that were

284 completely conserved across all species considered. Hence, the heatmap in Figure 4A is formed of the

285 2,782 proteins for which there were differences between two or more orthologous proteins in the

286 conservation of their identified phosphosites.

287 The dendrogram suggests three main groups for the *Plasmodium* species on the x-axis, with PPRFG01

288 (*Plasmodium sp.* gorilla clade 1) closest to the *P. falciparum* reference PF3D7 (Figure 4A). This group

289 (group 1) containing *P. falciparum* has six more species in the cluster, none of which (beyond 3D7) are

290 transmissible to humans. The two other groups are formed of 4 and 12 species. All other species

291 transmissible to humans are included in the largest group (group 3). On the y-axis, five clusters of

292 proteins were formed and subsequently analysed separately with clusterProfiler for enrichment

293 analysis, to determine if there were different biological processes associated with phosphosites of

294 different conservation patterns, which might be indicative of those under different selective

295 pressures. The analysis yielded 16 statistically significant GO terms (Figure 4B).

296 For Cluster 5 (proteins containing phosphosites mostly conserved across the genus), the most

297 significant GO terms per cluster were: "food vacuole" (GO:0020020), "extracellular vesicle"

298 (GO:1903561) and "endoplasmic reticulum" (ER) (GO:0005783) – indicating conserved signalling

299 mechanisms related to the infection of host cells. The ER in Apicomplexa is known to be involved in

300 the processing of effector proteins before translocation to host cells [31]. The "food vacuole" is a GO

301 term mostly used in annotation of Apicomplexa proteins related to digestion of the host cell

302 cytoplasm. Cluster 4 contains sites that are highly conserved in Group 1 species, and around 50%

303 conserved in Group 2 and 3 species. Cluster 4 is annotated to be enriched for GO terms related to RNA

304 binding and splicing. This result is somewhat surprising, since one would assume that mechanisms

305 related to transcription would be very highly conserved. Cluster 3 are highly conserved in Group 1

306 species but have low conservation in groups 2 and 3, with enrichment for GO terms related to

307 symbiont-containing vacuole membrane, rhoptry and extracellular vesicle – we could speculate that

308 this may be indicative of cell signalling related to host-cell invasion, and evolving faster that Cluster 5

309 proteins, to evade host immune responses. Cluster 2 proteins have highly conserved phosphosites in

310 Group 1 proteins but with average conservation between those in Cluster 2 and 4 for other groups.

311 Only one GO term is enriched for "response to xenobiotic stimulus" (GO:0009410) – a term related to

312 *Plasmodium*'s ability to respond to small molecules from the host (and proteins implicated in drug

313 resistance). Cluster 1 contains proteins with least conserved phosphosites, and had the strongest

314 enrichment (by p-value and count of mapped terms) for "infected host cell surface knob"

315 (GO:0020030) and "cell-cell adhesion" (GO:0098609). Term GO:0020030 is a commonly used gene

316 annotation in *Plasmodium*, related to the protrusions in the membrane of an infected erythrocyte.

317 These proteins are potentially under pressure to evade host immune responses, and thus evolving

318   much faster (and changing their cell signalling mechanisms). Data for this analysis including the

319   clusters can be found in Supp. File 4.

320

321   Figure 4. A: Heatmap of the agreement in sequence conservation between 23 species of *Plasmodium*: *P. gorilla* clade

322   G1(PPRFG01), *P. reichenowi* (PRCDC), *P. blacklocki* G01 (PBLACG01), *P. billcollinsi* G01 (PBILCG01), *P. adleri* G01 (PADL01), *P.*

323   *gaboni* (PGSY75), *P. malariae* (PmUG01), *P. brasilianum* strain Bolivian I (MKS88), *P. ovale* (PocGH01), *P. relictum* (PRELSG),

324   *P. gallinaceum* (PGAL8A), P. chabaudi AS (PCHAS), *P. vinckei vinckei* CY (PVVCY), *P. yoelii* 17X (PY17X), *P. berghei* ANKA

325   (PBANKA), *P. cynomolgi* (PcyM), *P. vivax* (PVP01), *P. knowlesi* (PKNH), *P. coatneyi Hackeri* (PCOAH), *P. fragile* strain nilgiri

326   (AK88), *P. inui* San Antonio 1 (C922), *P. vivax*-like (PVL) and the reference *P. falciparum* (PF3D7). The heatmap displays the

327   mean conservation (agreement) of phosphosites per protein with three clusters for species and five for proteins. B: Pathway

328   enrichment analysis for the genes found in each cluster with significant GO terms in y-axis and clusters in the x-axis; the

329   number of genes where those terms were found are in parenthesis under each cluster label.

330

331   We next implemented a "strict" phosphosite matching process, for the purposes of providing highly

332   likely phosphosites for all *Plasmodium* species aligned, requiring that the phosphosite amino acid

333   matched between *Plasmodium falciparum* 3D7 and the target species (allowing for S <-> T

334   substitutions) and requiring the +1 residue was also matched (as the most important position for

335   phosphorylation motifs). This gives an additional candidate set of 183,134 phosphosites, identified

336   across the *Plasmodium* genus, based on orthologue mapping (Supp. File 5). The multiple sequence

337   alignments for all phosphoproteins are provided in compressed folder (Supp. File 6).

338   We also investigated conservation based on SNP analysis within the *Plasmodium falciparum* species,

339   using data sets derived from whole genome sequencing of different isolates downloaded from

340   PlasmoDB (Supp. File 7). Out of 26,006 phosphosites, 25,587 did not have any recorded single amino

341   acid variants (SAAVs recorded in PlasmoDB), indicating very high conservation of phosphosites within

342   the species sampled (98.3%). Analysis of the total serines within the proteome found 4,660 SAAVs

343   with major allele frequency (AF) <1 (from 261,791 total serines), i.e. 98.2% conservation. This indicates

15

344     that that pS is no more likely or less likely to be mutated than other serines. On average, 97.9%

345     threonines in the proteome are conserved (i.e. have no SAAV in this analysis), compared to 98.7% for

346     pT sites. On average, 99.0% of pY (663/670) and 99.0% of all tyrosines do not have a SAAV in this

347     analysis, again indicating no particular selective pressure signal that could be identified. A histogram

348     of the major allele frequencies for phosphosite SAAVs is presented in Figure 5B, confirming that most

349     phosphosites are highly conserved across different isolates, with only a single site (pSer 33 in

350     PF3D7_1366900, a protein of unknown function) having major AF < 0.5. A table of proteins with

351     phosphosites and AF < 0.9 is provided in Supp. File 7, including several zinc finger proteins and two

352     rhoptry proteins.

353     We performed GO term enrichment analysis for the (419) proteins containing not fully conserved

354     phosphosites (Figure 5A). Besides "protein binding" (GO:0005515) or "mRNA binding" (GO:0003723)

355     which appear to be significant in most analyses, the analysis also returned "translocation of peptides

356     or proteins into host" (GO:0042000), "rhoptry" (GO:0020008) and "endocytosis" (GO:0006897),

357     indicative perhaps of proteins under some positive selective.

358     Other comparative analyses of conservation data have been included as Supp. Figure 2. In Supp Figure

359     2 (A) a pathway analysis was performed comparing three human transmissible species (PmUG01,

360     PocGH01, PVP01) vs. 17 not transmissible. *P. knowlesi* (PKNH) and *P. vivax*-like (PVL) were excluded

361     from this analysis because of their potential for non-vectorial infection to humans. From

362     phosphoproteins, a subset was selected for pathway analysis as: proteins with fully conserved

363     phosphosites for the human transmissible species, i.e. all phosphosites conserved with respect of the

364     reference PF3D7, and not conserved for non-transmissible species, at least one phosphosite within

365     the protein not conserved with respect to the *P. falciparum* reference PF3D7. Some of the significant

366     terms not previously observed in other analyses were "cytosol", "structural constituent of ribosome",

367     "cytosolic small ribosomal subunit". Supp Figure 2 (B) investigates the enriched pathways for sites

368     within phophoproteins not conserved in *P. gorilla* clade G1 (PPRFG01), compared to PF3D7, as our

369   analysis (Figure 4A) suggested this species to be the closest relative to *P. falciparium*. This additional

370   analysis returned "infected host cell surface knob" (GO:0020030) as the most significant term, which

371   may indicate differences on the cell invasion between the two species of *Plasmodium*.

372

373   Figure 5. A: Pathway analysis results across all *Plasmodium* species for those genes which were not fully conserved compared

374   to PF3D7, based on SNP analysis. B: Histogram of conservation across *Plasmodium* species based on SNP analysis.

### Disorder and structural analysis

376   Next, we explored the structural context of phosphosites in *P. falciparum*, to search for information

377   about the functional importance of the phosphosites. First, we performed an analysis to predict all

378   the structured and disordered regions of *P. falciparum* proteins (using metapredict v2), and mapped

379   phosphosites onto these regions. Metapredict gives a score from 0-1 to indicate the likelihood of each

380   amino acid within a sequence to be in an ordered (score = 0) or within a disordered region (score = 1).

381   It has been reported before that a high proportion of mammalian phosphosites are located on

382   disordered regions, having a role in transition from disorder to order, altering the local or global

383   structure, and potentially changing the interaction potential of the protein [32]. From our mapping of

384   phosphosites to predicted disordered regions, we could observe that phosphosites had a very strong

385   tendency towards disordered regions in *P. falciparum*. This tendency for phosphosites to be found in

386   disordered regions can be observed for all three residues S, T, Y. Figure 6A, displays boxplots with

387   strikingly higher disorder scores for phosphosites (pS, pT, pY) than for other S, T or Y residues in the *P.*

388   *falciparum* proteome – median disorder scores pS=0.982, pT = 0.965, and pY=0.975 compared to S=

389   0.769, T= 0.639 and Y= 0.532 median disorder scores of all residues in the proteome. The trend is

390   particularly striking for pY sites, since tyrosines do not have a strong preference to be in disordered

391   regions of proteins, as shown on Figure 6A. Metapredict documentation suggests that a threshold of

392   0.3 can differentiate ordered from disordered regions. Using such a threshold across the entire

393   proteome of *P. falciparum* would suggest that 66% of all residues fall in disordered regions. Comparing

394   against metapredict disorder scores for the human proteome (Supplementary Figure 3), revealed that

395    only 39% of residues within the human proteome are located in disordered regions. It is possible that

396    *P. falciparum* proteins are generally more disordered than human proteins, or that the tool is less well

397    calibrated for Apicomplexa than for humans. Nevertheless, using the metapredict-recommended

398    threshold of >0.3 for determining disordered regions, in agreement with previous reports [26],

399    revealed that 89% of phosphosites were located within predicted disordered regions (and still 85% if

400    a more conservative score >0.5 was used to determine disorder).

401    If we restrict the analysis to FLR "Gold" category phosphosites (those observed with high confidence

402    in more than one study), remarkably only 5.8% fall in ordered regions (on only 330 proteins), indicating

403    that it is highly unusual for phosphorylation to occur on well-structured regions of *P. falciparum*

404    proteins. For comparison, 12% of "gold standard" human phosphosites [23] are predicted to fall into

405    ordered regions (metapredict score < 0.3). Investigating the biological functions links to these 330

406    *Plasmodium falciparum* proteins with phosphosites in their ordered regions, there several GO terms

407    returned from pathway enrichment analysis with clusterProfiler (Supplementary Figure 4), including

408    proteins localising to the cytosol and cytoplasm, and ribosome-related functions.

409    Next, we wished to explore whether there was any difference in the conservation of phosphosites in

410    disordered vs ordered regions. As shown on Figure 6B, there is a striking difference – phosphosites in

411    ordered regions (2,776), had high conservation overall.

412    Examining the set of ordered sites, since these are relatively unusual in the *P. falciparum*

413    phosphoproteome, these sites are highly conserved (Figure 6B), compared to disordered phosphosites

414    – 48% of all "ordered" sites have conservation >90% across the genus, compared to only 16% of

415    "disordered" sites.

416    Exploring the high-quality (Gold) set of sites mapped to ordered regions (330 proteins), 209/330 (63%)

417    have a human ortholog (OrthoMCL DB [33]), indicative of genes highly conserved across all eukaryotes.

418    For the proteins containing "Gold" quality phosphosites in disordered regions (1,724), 718/1,724

419    (42%) have a human ortholog – indicative of proteins that are less well conserved across eukaryotes.

420    Disorder and site conservation data is available as Supp. File 8.

421    Figure 6. A: Boxplot of the disorder scores (from metapredict) for phosphosites (pS, pT and pY) versus disorder scores for all

422    S, T and Y residues in the *P. falciparum* proteome. B: Density functions for percentage of conservation by residue for ordered

423    and disordered regions (<0.3 metapredict score).

424

425    **Protein structural context**

426    The release of AlphaFold2 (AF2) has had a very significant impact on the ability to understand the

427    three-dimensional (3D) structure of proteins across both model and non-model organisms [34]. 3D

428    structure predictions for *P. falciparum* proteins are available via AlphaFold database, UniProtKB and

429    PlasmoDB (via mapping to UniProt identifiers). We have mapped all the phosphosites onto AF2

430    structures, also incorporating the conservation scores (across the *Plasmodium* genus), enabling visual

431    exploration of the relationship between order/disorder (which can be clearly visualized) on structures,

432    conservation and positions of phosphosites.  In Supp. File 9, we have created a static html page with

433    a table of the phosphoproteins identified, with hyperlinks so that every phosphoprotein's structure

434    and phosphosites can be visualized via the online iCn3D viewer [35], and links to their corresponding

435    record in UniProtKB (see section below on data access). A caveat is that AF2 models have no awareness

436    of PTM sites, and generally have been trained on few example proteins with PTMs intact. As such,

437    protein structure models demonstrate the position of PTMs as they would appear on an otherwise

438    unmodified structure. Given that phosphosites often change the structure of proteins by introducing

439    more negative charge, it remains an open research question how to re-model AF2 structures to reflect

440    the presence of phosphosites.

441    An example is presented in Figure 7, for protein phosphatase PPM2 (UniProtKB identifier: Q8IHY0,

442    PlasmoDB identifier: PF3D7_1138500). The image displays all identified phosphosites on the green to

443    red colour scale, indicating fully conserved across the genus = green, to unique to *P. falciparum* = red.

444    In PPM2, it can also be observed that phosphosites that are highly conserved (green) are located in

445    the structured core, and that large, disordered regions are around the outside, containing non-

446    conserved phosphosites in red.

447    When exploring functional relationships across orthologues, it might be typical to conclude that highly

448    conserved regions are most significant for function, and in fact most protein domains are captured

449    from conserved regions in multiple sequence alignments across assumed orthologues. However, the

450    data presented here does not support such a conclusion for phosphosites. The fact that the vast

451    majority of phosphosites are present on disordered regions of proteins, which are evolving fastest,

452    and seem to have a role in host cell invasion and potentially evading host cell responses, point to a

453    specialisation in the functional role of phosphorylation. Further experiments to understand the

454    interplay between phosphorylation, protein disorder and the ability to infect hosts are clearly

455    required.

456

457    Figure 7. Protein PF3D7_1138500 protein phosphatase PPM2, visualised in iCn3D (AlphaFold structure Q8IHY0), with mapped

458    phosphorylation sites – coloured red-black-green conservation scale (0-9). Phosphosites in the structured core are fully

459    conserved across the *Plasmodium* genus, whereas disordered regions have mostly low conservation.

460

461    **Open access data availability**

462    The *Plasmodium falciparum* "phosphosite build" is part of a wider project, called PTMeXchange, aimed

463    at high quality re-analysis of MS/MS data sets enriched for particular PTMs, and providing simple

464    public access to the resulting data sets. The build is provided via PRIDE/ProteomeXchange with

465    identifier PXD046874, where researchers can download the phosphosites identified per study in

466    simple tab-separated text files. The data has also been loaded into UniProtKB (Figure 8A), enabling

467    phosphosite evidence to be explored alongside other protein features and AF2 models, with links to

468    the raw evidence. We have also released the phosphosite build within PeptideAtlas

469    (https://peptideatlas.org/builds/pfalciparum/phospho/), enabling more detailed exploration of the

470    evidence within each protein, peptidoform and spectrum for a given site (Figure 8B). Our data are

471     deposited in PRIDE, UniProtKB and PeptideAtlas all provide evidence for each site using Universal

472     Spectrum Identifiers [36], which can be rendered via

473     https://proteomecentral.proteomexchange.org/usi/ (and shown in right panel on Figure 8B). This

474     allows a user to explore the evidence for a given site on a peptide, and test other possible explanations

475     to see if the spectrum could support alternative explanations (peptides or sites within those peptides).

476     Lastly, the phosphosite build is scheduled to appear in PlasmoDB in 2024, as this is the central

477     database for *Plasmodium* researchers.

478     Figure 8. A. Visualisation of phosphosites in UniProtKB for example protein (Q8IHY0). B. Examples of different visualisations

479     at the protein-, peptidoforms-, and spectral-levels for the same protein, with PeptideAtlas identifier PF3D7_1138500.1-p1.

480

**Discussion**

We have reported a comprehensive meta-analysis of *P. falciparum* phosphoproteomics data sets. Data sets were re-analysed using a robust pipeline ensuring objective false localisation rate calculation. Additional Gold, Silver, Bronze labelling was also provided as easy way to display confidence of phosphosites being correctly identified.

Results from our meta-analysis have been deposited into UniProtKB, enabling sites to be used for further research with other bioinformatics tools. The data has also made available in PeptideAtlas and PRIDE, enabling detailed exploration of scores and visualization of source mass spectra, as a full evidence trail.

We have also provided over-represented motifs found in these phosphoproteins as well as conservation data in relation to another 22 species of the genus *Plasmodium*, and from 115 different *Plasmodium falciparum* isolates, with respect to *P. falciparum* strain 3D7. We provided predictive disorder scores for all identified phosphosites and links to protein structures to visualise these proteins. We expect our results to be a useful resource for researchers facilitating identification of areas for future research in developing malaria treatments and vaccines.

**Methods**

**Data selection**

The ProteomeXchange Consortium [37] was used to identify suitable *P. falciparum* phosphoproteomics datasets, via the PRIDE repository [38] and ProteomeXchange. Overall 11 were deemed suitable for phosphosite (localisation) reanalysis, based on inclusion criteria that data sets were generated by data dependent acquisition (DDA) methods, "high-quality" (i.e. likely to deliver >1,000 phosphosites), raw data were available and readable using open source tools: PXD000070 [39], PXD001684 [4], PXD002266 [40], PXD005207 [41], PXD009157 [42], PXD009465 [43], PXD012143 [44], PXD015093 [45], PXD015833 [30], PXD020381 [46], and PXD026474 [47]. These 11 data sets consist of both labelled (TMT or iTRAQ) and unlabelled MS data sets. Most of the studies focus on the blood stage in *P. falciparum* life cycle, with only one study (PXD026474) including gametocytes. There was

508  not data available for the liver stage of the parasite, possibly due to challenges with the development

509  of liver models. A brief description of the studies' objectives extracted from the abstracts of their

510  publications, *P. falciparum* stages included in the analysis, enrichment details and MS methods are

511  shown in Table 2 and more details about the studies' objectives and biological samples in

512  Supplementary Table 1.

513

514    *Table 2. Summary results and sample used for each study included in this analysis.*

| Study ID | Summary | Sample | Ref |
|---|---|---|---|
| PXD000070 | "We analysed the Plasmodium *falciparum* schizont phosphoproteome using for the first time, a data-dependent neutral loss-triggered-ETD (DDNL) strategy and a conventional decision-tree method. … combination of Mascot Percolator and turbo-SLoMo represents a robust workflow for data analysis using CID and ETD fragmentation." | *P. falciparum* 3D7 was cultured in 2.5–5% O+ human erythrocytes. Infected erythrocytes were processed for analysis. | [39] |
| PXD001684 | "phosphoproteome analysis of extracellular merozoites revealing 1765 unique phosphorylation sites including 785 sites not previously detected in schizonts. " | *P. falciparum* 3D7 merozoites cultured *in vitro* | [4] |
| PXD002266 | "employing chemical and genetic tools in combination with quantitative global phosphoproteomics, we identify the phosphorylation sites on 69 proteins that are direct or indirect cellular targets for PfPKG. These PfPKG targets include proteins involved in cell signalling, proteolysis, gene regulation, protein export and ion and protein transport, indicating that cGMP/PfPKG acts as a signalling hub that plays a central role in a number of core parasite processes." | *P. falciparum* blood stage 3D7 (wild type)-, PKGT618Q- and CDPK1-HA- parasites were cultured *in vitro*, developing trophozoite and schizont stage parasites. For the time-course experiments, parasites were synchronized by two rounds of sorbitol treatment—first treatment when the parasites culture was at late ring/trophozoite stage and second when the parasite culture contained schizonts and ring stage parasites. | [40] |
| PXD005207 | "PfCDPK1 is critical for asexual development of *Plasmodium falciparum*, ... this kinase is critical for the invasion of host erythrocytes. Furthermore, using a multidisciplinary approach involving comparative phosphoproteomics we gain insights into the underlying molecular mechanisms." | *P. falciparum* 3D7 asexual blood stages were cultured in O+ erythrocytes. *P. falciparum* CDPK1- 3HA-DD parasites were cultured in the presence of 2.5 nM WR99210 and 0.25 µM Shld-1. | [41] |
| PXD009157 | "*Plasmodium falciparum* phosphodiesterase β (PDEβ) hydrolyses both cAMP and cGMP and is essential for blood stage viability. Conditional gene disruption causes a profound reduction in invasion of erythrocytes and rapid death of those merozoites that invade." | *P. falciparum* erythrocytic stages were cultured in human A+ erythrocytes | [42] |
| PXD009465 | "To better understand PfPK7-regulated phosphorylation events, we performed isobaric tag-based quantitative comparative phosphoproteomics of the schizont and segmenter stages from wild-type and pfpk7- parasite lines." | Both the *P. falciparum* parasite pfpk7– line and its parental 3D7 clone were grown in RPMI 1640 culture medium supplemented with A+ erythrocytes and 0.5% Albumax | [43] |
| PXD012143 | "investigated the role of cAMP in asexual blood stage development of Plasmodium falciparum through conditional disruption of adenylyl cyclase beta (ACβ) and its downstream effector, cAMP-dependent protein kinase (PKA)." | Both the *P. falciparum* parasite pfpk7– line and its parental 3D7 clone were grown in RPMI 1640 culture medium supplemented with A+ erythrocytes and 0.5% Albumax | [44] |
| PXD015093 | "global phosphoproteomic analysis of merozoites to identify signaling pathways that are activated during invasion." | *P. falciparum* 3D7 blood stages were cultured *in vitro* to generate merozoites | [45] |
| PXD015833 | Reports "species-specific phosphorylation of erythrocyte proteins by *P. falciparum*, but not by *Plasmodium knowlesi*, which does not export FIKK kinases." | *P. falciparum* parasites were cultured in complete media and *P. knowlesi* parasites were cultured in CM supplemented with 10% human serum | [30] |
| PXD020381 | "identify a multipass membrane protein, ICM1, with homology to transporters and calcium channels that is tightly associated with PKG in both asexual blood stages and transmission stages. Phosphoproteomic analyses reveal multiple ICM1 phosphorylation events dependent on PKG activity." | *P. falciparum* lines were maintained in human RBCs in RPMI 1640 containing AlbuMAX II | [46] |
| PXD026474 | "understand - CDPKs - role in human parasite transmission from the host to the mosquito vector and thus investigated the role of the human-infective parasite Plasmodium falciparum CDPK4 in the parasite life cycle. *P. falciparum* cdpk42 parasites created by targeted gene deletion showed no effect in blood stage development or gametocyte development. However, cdpk42 parasites showed a severe defect in male gametogenesis and the emergence of flagellated male gametes." | *P. falciparum* NF54 and Pfcdpk42 parasites were cultured as asexual blood stages according to standard procedures and received complete RPMI medium. Gametocytes were also cultured | [47] |

515

516    **Phosphosite localisation**

517    The search database was created from sequences derived from PlasmoDB 51 release of *Plasmodium*

518    *falciparum* 3D7 and human. The human sequences were obtained from the Level 1 PeptideAtlas Tiered

519 Human Integrated Search Proteome [48], containing core isoforms from neXtProt [49] (2020 build). A

520 fasta file was created combining these sequences with cRAP contaminant sequences

521 (https://www.thegpm.org/crap/) plus decoy sequences. One decoy sequence was generated each

522 protein and contaminant sequence using the de Bruijn method (with k=2) [50].

523 Data analysis was performed using the Trans-Proteomic Pipeline (TPP) [19], including Comet [51]

524 search engine for individual datasets, followed by PeptideProphet [52], iProphet [53], and

525 PTMProphet [54]; these were grouped within each study according to the labelling and MS

526 characteristics, as they determine search parameters. The files were searched with several variable

527 modifications depending on the study design, and phosphorylation of ASTY. Alanine was included in

528 searches as a decoy amino acid to enable the estimation of the FLR as explained by Ramsbottom *et al*.

529 [21]. Other variable modifications included in the analyses were: Oxidation in MWH (HYDR), protein

530 N-terminal acetylation or at K (ACET), ammonia loss at peptide n-terminal QC (PYRO), pyro-glu at

531 peptide n-terminal E (DHB), deamidation at NQ (DEAM), formylation of the N-terminus (FORM) and

532 fixed modification Carbamidomethylation (C), as noted in Supplementary Table 1. iTRAQ4plex, TMT6

533 or TMT10 labelling were included in searches as appropriate. As search parameters, a maximum of 2

534 missed cleavages specified and maximum number of 5 different modifications per peptide were

535 allowed.

536 **Post-processing search results**

537 The data files obtained from searching with TPP were processed by custom Python scripts

538 (https://github.com/PGB-LIV/mzidFLR) and analysed following a previously published pipeline [22].

539 First, at the peptide-spectrum match (PSM) level, the FDR was calculated based on decoy sequences

540 and the PSMs were filtered to 1% FDR. Data from most confident PSMs (1% FDR) were transformed

541 to give a site localisation score for each phosphosite found on each PSM, removing matches to decoy

542 PSMs and contaminant entries. A final site-based PSM score was obtained by multiplying the peptide

543 identification probability by the site localisation probability and adjusted considering the number of

544 occasions each site was observed, phosphorylated and not phosphorylated, in the dataset [22].

545     Once the final probability at PSM-level was calculated, the data was collapsed to "peptidoform-site"

546     level by taking the maximum score among all matches for each phosphosite on a given peptidoform

547     within each analysis and study. FLR was calculated by ordering all peptidoform-sites by their final

548     score. At this stage, matches to *P. falciparum* and human matches were separated according to their

549     respective protein matches.

550     A further categorisation at the protein-site level was achieved by using FLR calculations for each study.

551     This Gold-Silver-Bronze (GSB) categorisation allows additional grading for our confidence in the

552     findings among the most confident sites. A simple exclusive criterion was applied: Gold - sites

553     observed at 2 or more independent studies at 1% FLR; Silver - sites observed in only 1 dataset at <1%

554     FLR; Bronze: all other sites at 5% FLR.

555     Note that within each study, if there was more than 1 peptidoform-site as evidence for a protein site,

556     the one with lowest FLR was used for GSB categorisation.

557     **Downstream analysis**

558     **Summarising outcomes**

559     A frequency table was produced for PSMs at 1% FDR level, PSM-sites at 1% FDR level, overall number

560     of phosphosites at PSM-level at 1% and 5% FLR. Human and *P. falciparum* number of sites were

561     separated at peptidoform-site level and protein-site level at 1% and 5% FLR.

562     GSB phosphosites counts were generated considering that some of the sequences, and therefore sites,

563     could map to different proteins or sites within their respective proteome. The number of phosphosites

564     in sequences mapping to a single site (unique) and to more than one site were also calculated, as well

565     as the number of phosphosites mapping to only either the *P. falciparum* or Human proteome, or

566     belonging to sequences that map to both proteomes.

567     All analyses and graphical output were obtained using the R programming language (version 4.2.1) or

568     above, via RStudio (2022.02.3 Build 492).

**Motif and pathway enrichment analysis**

All *P. falciparum* datasets were investigated using motif and pathway enrichment analysis. For this, we used all STY phosphosites at 5% FLR, combining all 11 studies and removing matches to the decoy amino acid Alanine. 15mer peptides centred on each phosphosite were generated to investigate motifs around these sites. They were compared against 15mer "background" sequences generated using any STY in these datasets, whether they were phosphorylated or not, at its centre (position 0). Statistically significant enriched motifs surrounding phosphosites were identified using the R package *rmotifx* [55]. Results were thresholded via p-value < 1e-9 for pS and  p< 1e-6 for pT and pY respectively; and a minimum of 20 sequences per motif.

The R package clusterProfiler [56] was used to carry out a pathway enrichment analysis considering those phosphoproteins containing specific enriched motifs. All other proteins in the search database were used as background for this analysis. A heatmap on the adjusted p-values of a subset of motifs with fold-change enrichment (versus the background) > 4 was produced representing motifs against GO (Gene Ontology) terms.

**Conservation analysis**

*P. falciparum* isolate 3D7 identified phosphosites were investigated regarding 22 species of the *Plasmodium* genus. The 22 species of *Plasmodium* were: *P. gorilla clade G1*(PPRFG01), *P. reichenowi* (PRCDC), *P blacklocki* G01 (PBLACG01), *P. billcollinsi G01* (PBILCG01), *P. adleri G01* (PADL01), *P. gaboni* (PGSY75), *P. malariae* (PmUG01), *P. brasilianum strain Bolivian I* (MKS88), *P. ovale* (PocGH01), *P. relictum* (PRELSG), *P. gallinaceum* (PGAL8A), *P. chabaudi AS* (PCHAS), *P. vinckei vinckei CY* (PVVCY), *P. yoelii 17X* (PY17X), *P. berghei ANKA* (PBANKA), *P. cynomolgi* (PcyM), *P. vivax* (PVP01), *P. knowlesi* (PKNH), *P. coatneyi Hackeri* (PCOAH), *P. fragile strain nilgiri* (AK88), *P.inui San Antonio 1* (C922), *P. vivax-like* (PVL). They were compared against *P. falciparum* 3D7 in terms of sequence conservation. Protein sequences were downloaded from PlasmoDB and mapped to the identified phosphosites using the "syntenic ortholog" mappings stored in PlasmoDB, generated by OrthoMCL [33], followed by protein-level multiple sequence alignment with muscle 5.1 running on Linux [57].

595 Matched residues in orthologs were labelled as 1 while not matching sites, due to having a different

596 amino acid in that position in the ortholog, a gap in the sequence, or simply not having an ortholog

597 for a protein were labelled as 0. For each species and phosphoprotein considered, a mean

598 conservation score was calculated for the proportion of phosphosites within each protein, which were

599 conserved within the orthologue from that particular species. Heatmaps were created based on the

600 mean conservation score. Protein clusters were identified and investigated for enrichment analysis

601 using clusterProfiler.

602 Further functional enrichment analyses were performed comparing conservation results from human

603 transmissible species *P. malariae*, *P. ovale and P. vivax* to the other 17 animal species *(*excluding *P.*

604 *knowlesi* and *P. vivax-like* due to their zoonotic character). The subset of proteins included in the

605 enrichment analysis were those proteins fully conserved for the three human transmissible species

606 and those not fully conserved across all the animal species.

607 Conservation of phosphosites within *Plasmodium falciparum* was calculated using single nucleotide

608 polymorphism (SNP) data stored in PlasmoDB (public search strategy:

609 https://plasmodb.org/plasmo/app/workspace/strategies/import/b4d2489952494797), using

610 variants from 115 aligned genomes from two unbiased SNP data set

611 https://plasmodb.org/plasmo/app/record/dataset/DS_9a7f849906, and

612 https://plasmodb.org/plasmo/app/record/dataset/DS_d1c8287de9 [58]. SNP sites were

613 downloaded and non-synonymous variants, altering the amino acid sequence, were matched to the

614 phosphosite positions in proteins (using Python 3.9 code). Site conservation was estimated using the

615 major allele frequency, calculated by PlasmoDB.

616 We next explored the average disorder of amino acid positions within proteins using metapredict v2

617 [59] running on Linux, comparing and correlating disorder scores with site classifications. To visualise

618 examples of phosphosites with different conservation values, we used iCn3D [35]. The list of *P.*

619 *falciparum* 3D7 kinases was generated by searching for proteins annotated with Interpro or Pfam

620 "protein                    kinase"                    domains                    in                    PlasmoDB:

621 https://plasmodb.org/plasmo/app/workspace/strategies/import/6a11331d6eea4d20).

622 Phosphosite data for *Plasmodium falciparum* 3D7 was loaded into UniProtKB, by taking all mapping

623 all peptides carrying GSB sites to proteins in UniProtKB *P. falciparum* 3D7 proteome

624 (https://www.UniProt.org/proteomes/UP000001450), assuming tryptic cleavage.

625

634 **Author Contributions**

635 O.M.C. performed database searching, data analysis and manuscript writing. K.A.R. performed search

636 database generation, assisted with data processing. A.P. supported MS data curation. Y.P.R. supported

637 PRIDE data upload and curation. J.F., and M.M. assisted with data loading into UniProtKB and E.B-B

638 visualisation in UniProtKB. E.W.D and Z.S. assisted with MS data searching and creation of the

639 PeptideAtlas build. J.A.V. assisted with PRIDE data loading and project supervision. A.R.J. coordinated

640 the research, assisted with data analysis and writing the manuscript.

641

642 **References**

643 1.      Singh, B. and C. Daneshvar, *Human infections and detection of Plasmodium knowlesi.* Clin

644         Microbiol Rev, 2013. **26**(2): p. 165-84.

645    2.    Sinka, M.E., et al., *A global map of dominant malaria vectors.* Parasit Vectors, 2012. **5**: p. 69.

646    3.    Phillips, M.A., et al., *Malaria.* Nat Rev Dis Primers, 2017. **3**: p. 17050.

647    4.    Lasonder, E., et al., *Extensive differential protein phosphorylation as intraerythrocytic*

648          *Plasmodium falciparum schizonts develop into extracellular invasive merozoites.* Proteomics,

649          2015. **15**(15): p. 2716-2729.

650    5.    Pease, B.N., et al., *Global Analysis of Protein Expression and Phosphorylation of Three Stages*

651          *of Plasmodium falciparum Intraerythrocytic Development.* Journal of Proteome Research,

652          2013. **12**(9): p. 4028-4045.

653    6.    mal, E.R.A.C.G.o.V., *A research agenda for malaria eradication: vaccines.* PLoS Med, 2011.

654          **8**(1): p. e1000398.

655    7.    Arun Kumar, K., et al., *The circumsporozoite protein is an immunodominant protective antigen*

656          *in irradiated sporozoites.* Nature, 2006. **444**(7121): p. 937-940.

657    8.    Bonam, S.R., et al., *Plasmodium falciparum Malaria Vaccines and Vaccine Adjuvants.* Vaccines,

658          2021. **9**(10): p. 1072.

659    9.    Laurens, M.B., *RTS,S/AS01 vaccine (Mosquirix): an overview.* Hum Vaccin Immunother, 2020.

660          **16**(3): p. 480-489.

661    10.   Meshnick, S.R., T.E. Taylor, and S. Kamchonwongpaisan, *Artemisinin and the antimalarial*

662          *endoperoxides: from herbal remedy to targeted chemotherapy.* Microbiol Rev, 1996. **60**(2): p.

663          301-15.

664    11.   Nosten, F. and N.J. White, *Artemisinin-based combination treatment of falciparum malaria.*

665          Am J Trop Med Hyg, 2007. **77**(6 Suppl): p. 181-92.

666    12.   Wicht, K.J., S. Mok, and D.A. Fidock, *Molecular Mechanisms of Drug Resistance in Plasmodium*

667          *falciparum Malaria.* Annu Rev Microbiol, 2020. **74**: p. 431-454.

668    13.   Aurrecoechea, C., et al., *PlasmoDB: a functional genomic database for malaria parasites.*

669          Nucleic Acids Res, 2009. **37**(Database issue): p. D539-43.

670   14.   The Plasmodium Genome Database, C., *PlasmoDB: An integrative database of the Plasmodium falciparum genome. Tools for accessing and analyzing finished and unfinished sequence data. The Plasmodium Genome Database Collaborative.* Nucleic Acids Res, 2001. **29**(1): p. 66-9.

674   15.   Amos, B., et al., *VEuPathDB: the eukaryotic pathogen, vector and host bioinformatics resource center.* Nucleic Acids Res, 2022. **50**(D1): p. D898-D911.

676   16.   Gardner, M.J., et al., *Genome sequence of the human malaria parasite Plasmodium falciparum.* Nature, 2002. **419**(6906): p. 498-511.

678   17.   UniProt, C., *UniProt: the universal protein knowledgebase in 2021.* Nucleic Acids Res, 2021. **49**(D1): p. D480-D489.

680   18.   Urban, J., *A review on recent trends in the phosphoproteomics workflow. From sample preparation to data analysis.* Anal Chim Acta, 2022. **1199**: p. 338857.

682   19.   Deutsch, E.W., et al., *Trans-Proteomic Pipeline, a standardized data processing pipeline for large-scale reproducible proteomics informatics.* Proteomics Clin Appl, 2015. **9**(7-8): p. 745-54.

685   20.   Keller, A., et al., *Empirical Statistical Model To Estimate the Accuracy of Peptide Identifications Made by MS/MS and Database Search.* Analytical Chemistry, 2002. **74**(20): p. 5383-5392.

687   21.   Ramsbottom, K.A., et al., *Method for Independent Estimation of the False Localization Rate for Phosphoproteomics.* J Proteome Res, 2022. **21**(7): p. 1603-1615.

689   22.   Camacho, O.M., et al., *Assessing multiple evidence streams to decide on confidence for identification of post-translational modifications, within and across data sets.* bioRxiv, 2022: p. 2022.12.15.520504.

692   23.   Kalyuzhnyy, A., et al., *Profiling the Human Phosphoproteome to Estimate the True Extent of Protein Phosphorylation.* Journal of Proteome Research, 2022. **21**(6): p. 1510-1524.

694   24.   Perez-Riverol, Y., et al., *The PRIDE database resources in 2022: a hub for mass spectrometry-based proteomics evidences.* Nucleic Acids Research, 2021. **50**(D1): p. D543-D552.

696  25.  Desiere, F., et al., *The PeptideAtlas project.* Nucleic Acids Res, 2006. **34**(Database issue): p.
697       D655-8.

698  26.  Treeck, M., et al., *The phosphoproteomes of Plasmodium falciparum and Toxoplasma gondii*
699       *reveal unusual adaptations within and beyond the parasites' boundaries.* Cell Host Microbe,
700       2011. **10**(4): p. 410-9.

701  27.  Adderley, J. and C. Doerig, *Comparative analysis of the kinomes of Plasmodium falciparum,*
702       *Plasmodium vivax and their host Homo sapiens.* BMC Genomics, 2022. **23**(1): p. 237.

703  28.  Johnson, J.L., et al., *An atlas of substrate specificities for the human serine/threonine kinome.*
704       Nature, 2023. **613**(7945): p. 759-766.

705  29.  Ward, P., et al., *Protein kinases of the human malaria parasite Plasmodium falciparum: the*
706       *kinome of a divergent eukaryote.* BMC Genomics, 2004. **5**: p. 79.

707  30.  Davies, H., et al., *An exported kinase family mediates species-specific erythrocyte remodelling*
708       *and virulence in human malaria.* Nat Microbiol, 2020. **5**(6): p. 848-863.

709  31.  Coffey, M.J., et al., *Role of the ER and Golgi in protein export by Apicomplexa.* Curr Opin Cell
710       Biol, 2016. **41**: p. 18-24.

711  32.  Newcombe, E.A., et al., *How phosphorylation impacts intrinsically disordered proteins and*
712       *their function.* Essays Biochem, 2022. **66**(7): p. 901-913.

713  33.  Chen, F., et al., *OrthoMCL-DB: querying a comprehensive multi-species collection of ortholog*
714       *groups.* Nucleic acids research, 2006. **34**(Database issue): p. D363-8.

715  34.  Jumper, J., et al., *Highly accurate protein structure prediction with AlphaFold.* Nature, 2021.
716       **596**(7873): p. 583-589.

717  35.  Wang, J., et al., *iCn3D, a web-based 3D viewer for sharing 1D/2D/3D representations of*
718       *biomolecular structures.* Bioinformatics, 2019. **36**(1): p. 131-135.

719  36.  Deutsch, E.W., et al., *Universal Spectrum Identifier for mass spectra.* Nature Methods, 2021.
720       **18**(7): p. 768-770.

721   37.   Vizcaino, J.A., et al., *ProteomeXchange provides globally coordinated proteomics data submission and dissemination.* Nat Biotechnol, 2014. **32**(3): p. 223-6.

723   38.   Martens, L., et al., *PRIDE: the proteomics identifications database.* Proteomics, 2005. **5**(13): p. 3537-45.

725   39.   Collins, M.O., et al., *Confident and sensitive phosphoproteomics using combinations of collision induced dissociation and electron transfer dissociation.* J Proteomics, 2014. **103**(100): p. 1-14.

728   40.   Alam, M.M., et al., *Phosphoproteomics reveals malaria parasite Protein Kinase G as a signalling hub regulating egress and invasion.* Nature Communications, 2015. **6**(1): p. 7285.

730   41.   Kumar, S., et al., *PfCDPK1 mediated signaling in erythrocytic stages of Plasmodium falciparum.* Nat Commun, 2017. **8**(1): p. 63.

732   42.   Flueck, C., et al., *Phosphodiesterase beta is the master regulator of cAMP signalling during malaria parasite invasion.* PLoS Biol, 2019. **17**(2): p. e3000154.

734   43.   Pease, B.N., et al., *Characterization of Plasmodium falciparum Atypical Kinase PfPK7(-) Dependent Phosphoproteome.* J Proteome Res, 2018. **17**(6): p. 2112-2123.

736   44.   Patel, A., et al., *Cyclic AMP signalling controls key components of malaria parasite host cell invasion machinery.* PLoS Biol, 2019. **17**(5): p. e3000264.

738   45.   More, K.R., et al., *Phosphorylation-Dependent Assembly of a 14-3-3 Mediated Signaling Complex during Red Blood Cell Invasion by Plasmodium falciparum Merozoites.* mBio, 2020. **11**(4).

741   46.   Balestra, A.C., et al., *Ca2+ signals critical for egress and gametogenesis in malaria parasites depend on a multipass membrane protein that interacts with PKG.* Science Advances, 2021. **7**(13): p. eabe5396.

744   47.   Kumar, S., et al., *Plasmodium falciparum Calcium-Dependent Protein Kinase 4 is Critical for Male Gametogenesis and Transmission to the Mosquito Vector.* mBio, 2021. **12**(6): p. e0257521.

747   48.   Deutsch, E.W., et al., *Tiered Human Integrated Sequence Search Databases for Shotgun*
748         *Proteomics.* J Proteome Res, 2016. **15**(11): p. 4091-4100.

749   49.   Lane, L., et al., *neXtProt: a knowledge platform for human proteins.* Nucleic Acids Res, 2012.
750         **40**(Database issue): p. D76-83.

751   50.   Moosa, J.M., et al., *Repeat-Preserving Decoy Database for False Discovery Rate Estimation in*
752         *Peptide Identification.* J Proteome Res, 2020. **19**(3): p. 1029-1036.

753   51.   Eng, J.K., T.A. Jahan, and M.R. Hoopmann, *Comet: an open-source MS/MS sequence database*
754         *search tool.* Proteomics, 2013. **13**(1): p. 22-4.

755   52.   Ma, K., O. Vitek, and A.I. Nesvizhskii, *A statistical model-building perspective to identification*
756         *of MS/MS spectra with PeptideProphet.* BMC Bioinformatics, 2012. **13**(16): p. S1.

757   53.   Shteynberg, D., et al., *iProphet: multi-level integrative analysis of shotgun proteomic data*
758         *improves peptide and protein identification rates and error estimates.* Mol Cell Proteomics,
759         2011. **10**(12): p. M111 007690.

760   54.   Shteynberg, D.D., et al., *PTMProphet: Fast and Accurate Mass Modification Localization for*
761         *the Trans-Proteomic Pipeline.* J Proteome Res, 2019. **18**(12): p. 4262-4272.

762   55.   Wagih, O., et al., *Uncovering Phosphorylation-Based Specificities through Functional*
763         *Interaction Networks*.* Molecular & Cellular Proteomics, 2016. **15**(1): p. 236-245.

764   56.   Yu, G., et al., *clusterProfiler: an R package for comparing biological themes among gene*
765         *clusters.* Omics, 2012. **16**(5): p. 284-7.

766   57.   Edgar, R.C., *MUSCLE: a multiple sequence alignment method with reduced time and space*
767         *complexity.* BMC Bioinformatics, 2004. **5**(1): p. 113.

768   58.   Chang, H.H., et al., *Malaria life cycle intensifies both natural selection and random genetic*
769         *drift.* Proc Natl Acad Sci U S A, 2013. **110**(50): p. 20129-34.

770   59.   Emenecker, R.J., D. Griffith, and A.S. Holehouse, *Metapredict: a fast, accurate, and easy-to-*
771         *use predictor of consensus disorder and structure.* Biophysical Journal, 2021. **120**(20): p. 4312-
772         4319.

773 **Supplemental Information**

774 Supp File 1. GSB_withMotifs.csv : Phosphosites at 5% FLR classified as Gold, Silver and Bronze and

775 statistically significant motifs for those sites.

776 Supp File 2. EnrichmentResult_Allmotifs_05.xlsx : Enrichment analysis results for all statistically

777 significant motifs.

778 Supp File 3. Conservation.csv : Conservation scores across species with respect to the reference 3D7.

779 Supp File 4. ConservationSpeciesWithClusters.csv : Proportion of sites conserved within each protein

780 and species with respect to 3D7 and its clusters.

781 Supp File 5. mapped_plasmodium_sites.tsv : Phosphosites putatively identified in other *Plasmodium*

782 species based on orthologue mapping.

783 Supp File 6. plasmodium_alignments.zip : Multiple sequence aligments of orthologous

784 phosphoproteins within the *Plasmodium* genus.

785 Supp File 7. SNP_data.xlsx : Conservation analysis based on single amino acid variants.

786 Supp File 8. Disorder.csv : Disorder scores from metapredict for Gold, Silver and Bronze phosphosites.

787 Supp File 9. Proteins_3D : Hyperlinks for visualising identified phosphoproteins in iCn3D viewer.

788 Supplementary Information.docx : Suplementary tables and figures.

## Group1: 57 Motifs

......[ST]D......
......D[ST]........
.........SE........
......ES........
......[ST].[DE].....
......S ..[DE].....

......T ..E....
...[DE]..S .......
......[ST]...D...
......S...[DE]...
.....D.TE......
....E..TE......

## Group2: 23 Motifs

......SN......
......N[ST]......
...K..S.N....
....N.S..[DE]....
....R..S..N....
.....N.SD......
.....N.SP......

## Group3: 15 Motifs

...K..[ST]......
......[ST].S.....
......GS........
.......S.S......

## Group4: 12 Motifs

.......Y.......

## Group5: 8 Motifs

......[ST]P......

## Group6: 11 Motifs

....R..[ST].......
......R.S.......

**Group 1**

**Group 2**

**Group 3**

**Group 4**

**Group 5**

**Group 6**

A

A

B