

1 **Pervasive adaptation in *Plasmodium*-interacting proteins in mammals**

2

3 Emily R. Ebel^{1,2}, Natalie Telis³, Sandeep Venkataram^{1,4}, Dmitri A. Petrov¹, David

4 Enard^{1,2}

5

6 ¹ Department of Biology, Stanford University, Stanford, California 94305

7 ² Corresponding Author

8 ³ Program in Biomedical Informatics, Stanford University, Stanford, California 94305

9 ⁴ Present Address: Division of Biological Sciences, University of California San Diego,

10 San Diego, California 97092

11

12 Contact:

13 Emily Ebel, David Enard

14 371 Serra Mall, Stanford, CA, 94305

15 (650) 736-2249

16 ebel@stanford.edu; denard@stanford.edu

17

18

19

20

21

22

1 **Abstract**

2 The protozoan genus *Plasmodium* causes malaria in dozens of mammal
3 species, including humans, non-human primates, rodents, and bats. In humans,
4 *Plasmodium* infections have caused hundreds of millions of documented deaths,
5 imposing strong selection on certain populations and driving the emergence of
6 several resistance alleles. Over the deep timescale of mammalian evolution,
7 however, little is known about host adaptation to *Plasmodium*. In this work, we
8 expand the collection of known *Plasmodium*-interacting-proteins (PIPs) in
9 mammalian hosts from ~10 to 410, by manually curating thousands of scientific
10 abstracts. We use comparative tests of adaptation to show that PIPs have
11 experienced >3 times more positive selection than similar mammalian proteins,
12 consistent with *Plasmodium* as a major and long-standing selective pressure. PIP
13 adaptation is strongly linked to gene expression in the blood, liver, and lung, all of
14 which are clinically relevant tissues in *Plasmodium* infection. Interestingly, we find
15 that PIPs with immune functions are especially enriched for additional interactions
16 with viruses or bacteria, which together drive a 3.7-fold excess of adaptation. These
17 pleiotropic interactions with unrelated pathogens, along with pressure from other
18 *Plasmodium*-like Apicomplexan parasites, may help explain the PIP adaptation we
19 observe in all clades of the mammalian tree. As a case study, we also show that
20 alpha-spectrin, the major membrane component of mammalian red blood cells, has
21 experienced accelerated adaptation in domains known to interact specifically with
22 *Plasmodium* proteins. Similar interactions with *Plasmodium*-like parasites appear to

- 1 have driven substantial adaptation in hundreds of host proteins throughout
- 2 mammalian evolution.

1 Introduction

2 Malaria is one of the world's most notorious infectious diseases, responsible
3 for billions of illnesses and millions of deaths in the last fifty years alone (WHO,
4 2015). The malaria genus *Plasmodium* contains five species infecting humans,
5 including *P. falciparum*, and 53 species infecting non-human primates, rodents, and
6 bats (Carlton, Perkins, and Deitsch, 2013). Other blood-borne Apicomplexans, which
7 are frequently confused for *Plasmodium*, also cause malaria-like symptoms in
8 livestock and pets (Escalante and Ayala, 1995; Coatney and Roudabush, 1936; Clark
9 and Jacobson, 1998). The strict genus *Plasmodium* is thought to have experienced a
10 major radiation 55-129 million years ago (Escalante and Ayala, 1995), indicating a
11 long-standing relationship between the malaria parasite and its mammalian hosts.

12 This ancient relationship raises the possibility that *Plasmodium*, along with
13 similar parasites, have imposed an important and complex selective pressure on
14 mammals. The host-parasite interactions involved in malaria span multiple stages
15 and tissue types, each of which may be subject to selection. Briefly, after the bite of
16 an infected mosquito transmits *Plasmodium* cells into the blood, they migrate to the
17 liver and multiply many times. After several days, parasites emerge from the liver
18 and begin infecting red blood cells (RBCs). The ensuing 48-hour cycles of replication
19 and emergence from RBCs are responsible for anemia, fever, and other
20 characteristic symptoms of malaria. *Plasmodium* parasites are also known to
21 sequester in certain organs, including the brain, lungs, and adipose tissues, which
22 can result in severe complications (Idro et al. 2010; Franke-Fayard et al., 2005;
23 Lovegrove et al., 2008; Aursudkij et al., 1998). Furthermore, parasitic proteins and

1 by-products solicit a complex immune response, including the tagging of parasitized
2 RBCs for removal from circulation by the spleen (Engwerda *et al.*, 2005).

3 Given these many facets of host-parasite interaction, as well as the
4 substantial morbidity and mortality of malaria, it seems likely that *Plasmodium* has
5 imposed an important selective pressure on its hosts. Over the deep time scale of
6 mammalian evolution, this hypothesis has not yet been tested, but it has been
7 supported over the shorter time scale of human evolution. In African and Southeast
8 Asian populations, several malaria resistance variants appear to have risen in
9 frequency over the last 5,000-10,000 years (Hedrick, 2011; Kwiatkowski, 2005).
10 While some, such as the Duffy null mutation, have approached local fixation (Welch,
11 1977), many others are prevented from fixing by their deleterious pleiotropic
12 effects. For example, the hemoglobin sickle cell allele offers substantial protection
13 against malaria, but causes fatal anemia in the homozygous form (Aidoo *et al.*,
14 2002). The fact that balancing selection maintains such a deleterious allele, at up to
15 15% frequency in some African populations (Piel *et al.*, 2010), suggests that malaria
16 presents a strongly opposing selective force. Indeed, malaria has repeatedly been
17 labeled "one of the strongest selective forces on the human genome" (Hedrick, 2011;
18 Verra *et al.*, 2009; Kwiatkowski, 2005), though this statement has never been
19 quantified.

20 The evolutionary impact of a complex selective pressure is difficult to
21 quantify precisely. One important reason is that the many genes relevant to one
22 phenotype, like malaria resistance, are each likely to be pleiotropically involved
23 with a set of other phenotypes (Wagner and Zhang, 2011). Thus, it is difficult to

1 ascribe evolutionary patterns in a small number of genes with malaria-related
2 functions, such as certain RBC or immune proteins, specifically to the selective
3 pressure imposed by malaria. Limiting evolutionary analysis to certain pathways or
4 genes can also exclude the effects that a complex selective pressure has on other
5 biological systems (see Trivisano and Shaw, 2012). Both of these issues could be
6 circumvented, in the case of malaria, by comparing a large set of genes known to
7 interact with *Plasmodium* to a large set of genes that share similar properties, but do
8 not interact with *Plasmodium*. With a sufficient number of genes, this approach
9 should average out various pleiotropic effects, as well as allow for expression in
10 particular tissues to be tested for association with adaptation. A similar approach to
11 gene set enrichment has previously been used to describe polygenic adaptation to
12 pathogens in humans (Daub et al. 2013). More recently, it has been combined with
13 manual curation and carefully chosen controls to identify viruses as a dominant
14 driver of adaptation in mammalian proteins (Enard et al., 2016).

15 In order to estimate an evolutionary effect specifically attributable to
16 *Plasmodium*, the gene enrichment strategy outlined above necessarily requires a
17 large number of *Plasmodium*-interacting genes. However, only a few dozen such
18 genes had previously been compiled (e.g. Verra et al., 2009), and only ~10 of these
19 are highly conserved across mammals. In this work, we manually examine over
20 30,000 publications to identify 410 conserved, mammalian *Plasmodium*-interacting
21 proteins, or PIPs. By leveraging extensive collections of functional data from over
22 9,000 mammalian proteins, we fairly compare rates of evolution between PIPs and
23 control genes, which have been matched to PIPs across a wide range of properties.

1 We find evidence of unusually strong and pervasive positive selection in PIPs, which
2 has likely been driven by interactions with *Plasmodium*-like parasites over millions
3 of years.

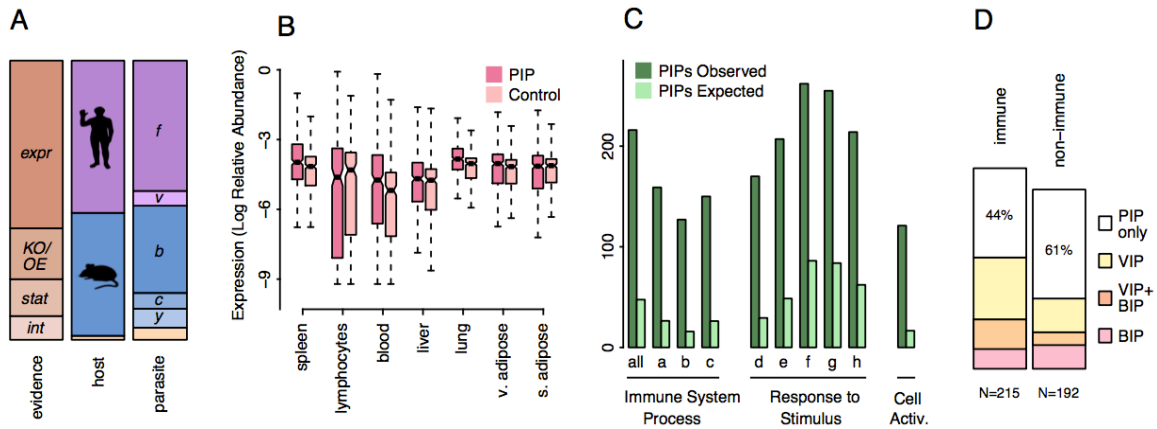
4

5 **Results**

6 ***Robust identification of 410 Plasmodium-interacting proteins (PIPs)***

7 Given the importance of malaria to human populations, research into host-
8 parasite interactions and host health effects has been underway for over a century.
9 We queried the PubMed database for all scientific papers whose abstracts
10 mentioned 'malaria,' '*falciparum*,' or '*Plasmodium*,' along with the name of a host
11 gene (Methods I). To focus on mammalian evolution, we limited our search to 9,338
12 protein-coding genes that are conserved in 24 mammalian species, including
13 humans (Methods II; Enard et al., 2016). Most of these mammalian species belong to
14 one of four orders—primates, rodents, artiodactyls, or carnivores (S1 Fig)—and
15 represent a variety of lifestyles and susceptibilities to malaria.

16 This search returned 30,788 papers associated with 2,249 of the 9,338 well-
17 conserved genes. However, the vast majority of these results were false positives,
18 due largely to multiple meanings of the short acronyms that identify genes. We
19 manually curated the results to identify 484 papers linking 410 proteins to malaria
20 through any of four types of phenotypic evidence: (1) physical interaction between
21 the mammalian and *Plasmodium* protein; (2) statistical association with malaria
22 susceptibility; (3) knock-out or overexpression studies; and (4) low-throughput
23



1

2 **Fig. 1. PIP sources, expression patterns, and GO functions.** (a) The 410 PIPs are summarized in
 3 each stacked bar. Abbreviations: *expr* – expression; *KO/OE* – knockout/overexpression; *stat* –
 4 statistical association; *int* - physical interaction; *f* - *falciparum*; *v* - *vivax*; *b* - *berghei*; *c* - *chabaudi*; *y* -
 5 *yoelii*. (b) Relative abundance of PIP expression is higher than expected in seven tissues (all $p < 0.006$),
 6 compared to control genes matched for total expression. All tissues are shown in S2 Fig.
 7 Abbreviations: *v* - visceral; *s* - subcutaneous. (c) Top 10 GO categories most significantly enriched for
 8 PIPs (all $p < 0.001$). Sub-categories a-h are enumerated in S3 Fig. (d) The pathogen pleiotropy of
 9 immune PIPs (53% of all PIPs) and non-immune PIPs (47% of all PIPs) is summarized in each
 10 stacked bar. 44% of immune PIPs, versus 61% of non-immune PIPs, are not known to interact with
 11 any other viral (VIP) or bacterial (BIP) pathogen.

12

13

14 studies showing a change in gene expression during malaria infection (Fig 1A; S1
 15 Table). Expression changes were the most common form of evidence (72% of PIPs),
 16 but 28% of PIPs were supported by multiple sources of evidence, and 41% by
 17 multiple studies. Virtually all of the studies were conducted on five *Plasmodium*
 18 species infecting humans or mice (Fig 1A).

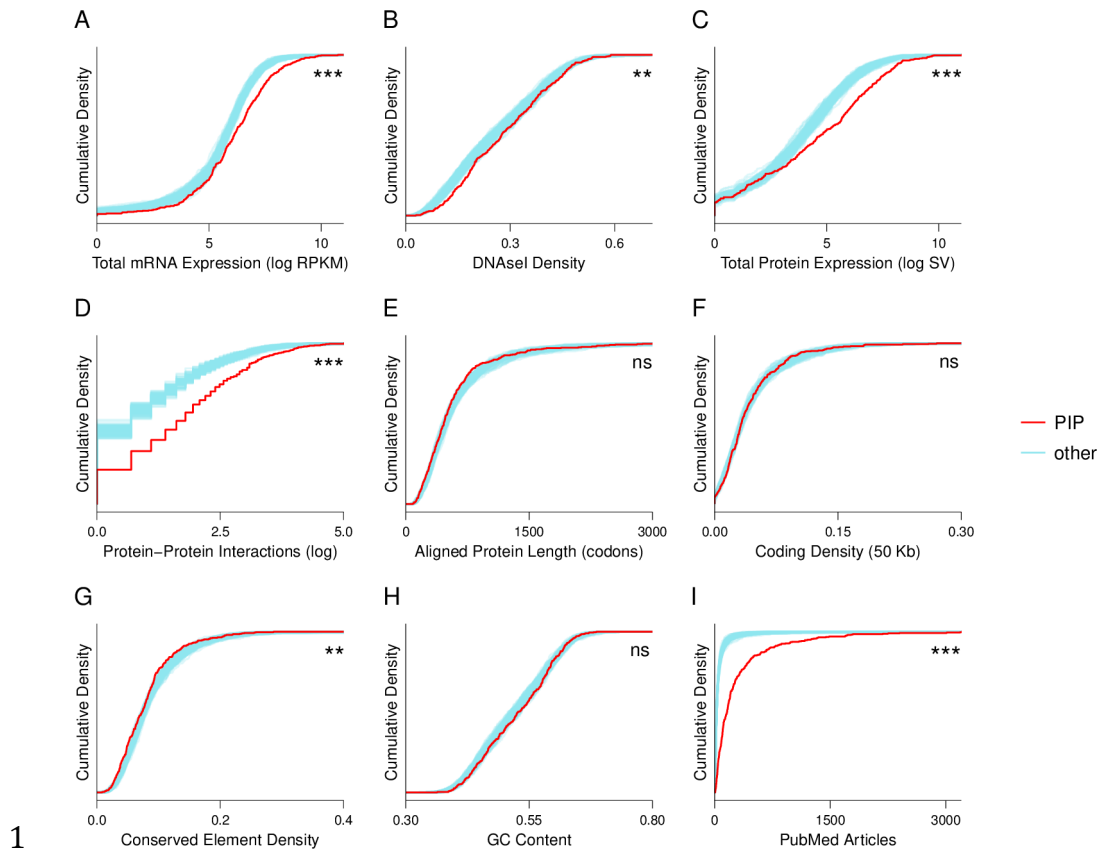
19

20

1 ***PIP expression and function support role in malaria***

2 If PIPs are truly a set of malaria-relevant genes, we would expect the
3 pathophysiology of malaria to be reflected in their tissue expression profiles. We
4 tested this hypothesis by examining human gene expression in each of the 53 tissues
5 collected by the GTEx Consortium (2015; Methods III). We first found that, on
6 average, PIPs have 9.5% higher total expression than other genes ($p < 0.0001$; Fig
7 2A). To fairly evaluate PIP overexpression in each tissue, we designed a matched
8 permutation test that compares PIPs to many, similarly-sized sets of control genes
9 with similar total expression (Methods III). After controlling for total expression in
10 this way, we find seven tissues in which PIPs are significantly overexpressed (Fig
11 1B, all $p < 0.006$; all tissues are shown in S2 Fig). These include blood and liver,
12 where *Plasmodium* parasites reproduce; spleen and lymphocytes, which participate
13 in the immune response; and the lung and adipose tissues, which are secondary
14 sites of parasite sequestration (Franke-Fayard et al., 2005).

15 Similarly, we expected PIPs to be enriched for GO functions that reflect
16 malaria pathology. We tested 17,696 GO functional categories (Methods V) for PIP
17 enrichment using Fisher's Exact Test. After correcting for multiple testing, over
18 1,000 categories contained significantly more PIPs than expected (S2 Table). These
19 categories are dominated by immune functions, especially for the highest levels of
20 enrichment (Fig 1C). Other functions, including apoptosis, cell-cell signaling, and
21 coagulation, are also highly enriched for PIPs (S2 Table). These results confirm the
22 biological connections between PIPs and malaria, and suggest that immune
23 pathways present a major functional interface between host and parasite.



1

2

Fig. 2. PIPs differ from other conserved genes by a number of functional measures, but meet

3

expectations for most measures of genomic context. Each plot is a cumulative density function, which

4

tracks the proportion of the gene set (y-axis) with values greater than or equal to each value on the x-

5

axis. That is, sets with higher average values will traverse more x-axis space (and appear as 'lower'

6

lines) before reaching the maximum density of 100%. Descriptions of data sources are available in

7

Methods V. *** = $p < 0.0001$; ** = $p < 0.001$; ns = $p > 0.05$.

8

9

10 ***Immune and other PIPs are enriched for viral and bacterial interactions***

11

Many immune genes, even outside the adaptive immune system, are known

12

to be activated by signals from multiple pathogens (e.g. Ozinsky et al, 2000;

13

Yamamoto et al., 2013). Such 'pathogen pleiotropy' poses an important complication

14

when testing the link between *Plasmodium*, as a single causal pathogen, and

1 adaptation in any gene. To quantify the extent of this pleiotropy for PIPs, we
2 compiled mammalian proteins known to interact with viruses, bacteria, and any
3 Apicomplexan parasites outside the genus *Plasmodium* (Methods IV). For viruses,
4 we obtained a high-quality list of 1,256 manually curated virus-interacting proteins
5 (VIPs) from Enard *et al* (2016). For bacteria and Apicomplexans, we queried the EBI
6 IntAct database (Orchard et al., 2014) for all deposited interactions (see Methods
7 IV). This search returned 1,250 mammalian bacteria-interacting proteins (BIPs; S3
8 Table), but 0 instances of mammalian interaction with non-*Plasmodium*
9 Apicomplexans. Most Apicomplexan parasites are still poorly studied at the
10 molecular level (Carlton, Perkins, and Deitsch, 2013), but it is likely that parasites in
11 the same phylum as *Plasmodium* have more overlapping interactions with PIPs than
12 do viruses or bacteria.

13 Overall, we find that 37% of all PIPs also interact with viruses, 22% with
14 bacteria, and 48% with viruses and/or bacteria—many more than expected by
15 chance (all $p < 0.0001$; S6 Fig). As expected, this multi-pathogen overlap is strongest
16 for immune PIPs (Fig 1D). Surprisingly, however, we find that nearly 40% of non-
17 immune PIPs also interact with these unrelated groups of pathogens (Fig 1D). While
18 some of these non-immune, multi-pathogen PIPs could in theory have
19 uncharacterized immune functions, most are known for their involvement in
20 general cellular processes, including metabolism and signal transduction. This
21 suggests that a diverse array of prokaryotic, eukaryotic, and viral pathogens take
22 advantage of a limited number of cellular pathways to infect their hosts. Such

1 pleiotropy has many interesting implications, including the need to carefully isolate
2 any single selective pressure when linking it to protein adaptation.

3

4 ***PIPs are not like other proteins***

5 We have already shown that PIPs have two unusual properties—high mRNA
6 expression, and excess overlap with other pathogens—that may influence their rate
7 of evolution. We assessed several additional metrics for differences between PIPs
8 and other proteins, in order to fairly evaluate PIP adaptation.

9 First, we tested three more broad measures of gene function in humans: the
10 density of DNaseI hypersensitive elements; protein expression, as measured by
11 mass spectrometry; and the number of protein-protein interactions (see Methods
12 V). For each of these metrics, PIPs have significantly higher mean values than sets of
13 random controls, indicating that PIPs are more broadly functional in humans (Fig 2,
14 B-D; all $p < 0.01$). We next tested four measures of genomic context, which have been
15 linked to the rate of protein evolution: aligned protein length; the regional density of
16 protein-coding bases; the density of highly conserved, vertebrate elements; and GC
17 content (Methods V). Most of these metrics do not differ between PIPs and other
18 genes (Fig 2 E-H), with the exception of conserved element density, which is slightly
19 but significantly lower in PIPs (mean=8.0% vs. 8.8%; $p=0.0004$; Fig 2G).

20 Based on these results, we expanded our permutation test to find matched
21 controls for each PIP. Control genes were considered acceptable matches if their
22 values for each of the five significantly different metrics (Fig 2A-D; G) fell within
23 specific ranges of the PIP value (Methods VI). This permutation procedure

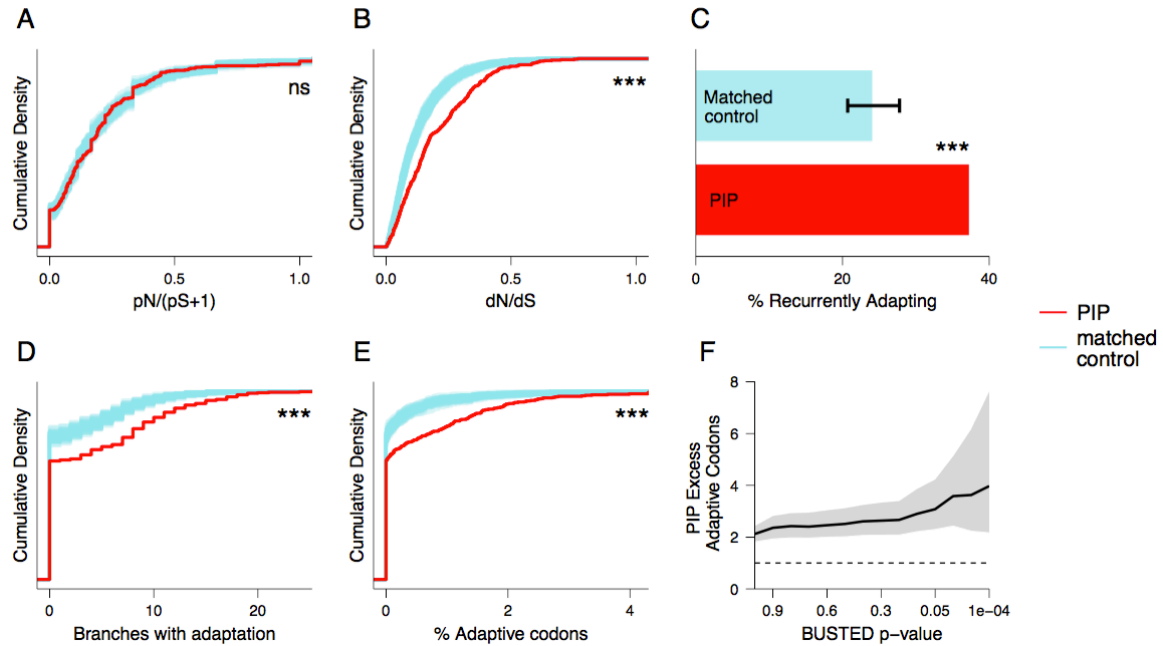
1 effectively equalized PIP and control distributions for all eight functional and
2 evolutionary properties displayed in Fig 2A-H (S4 Fig). On average, each PIP could
3 be matched to 32 control genes, allowing many different sets of matched controls to
4 be generated. About 9% of PIPs were too dissimilar from other proteins to be
5 matched, and were excluded from subsequent analysis.

6 Finally, one of the largest differences between PIPs and other proteins is the
7 frequency with which they are discussed in the scientific literature (Fig 2I). The
8 average PIP has 6.5 times more PubMed citations, and 9.1 times more scientist-
9 contributed References Into Function (GeneRIFs), than the average mammalian
10 protein (Methods V). This difference was too large to control for in the matched
11 permutation test without excluding the majority of PIPs. However, we show that the
12 citation frequency of non-PIPs has no relationship with measures of protein
13 adaptation ($p \geq 0.3$; S5 Fig, A&B). Furthermore, non-PIPs in the top quintile of
14 citation frequency have no more adaptation than other genes ($p \geq 0.25$; S5 Fig, C&D).
15 This indicates that a high rate of citation for PIPs is not likely to be causally
16 associated with their rate of adaptation.

17

18 ***PIPs have experienced accelerated rates of adaptation in mammals***

19 If malaria parasites (or similar pathogens) have indeed imposed a major
20 selective pressure on their mammalian hosts, we would expect PIPs to exhibit
21 unusual, adaptive patterns of amino acid substitution. In the absence of exceptional
22 selection, however, these patterns would not be expected. After controlling for



1

2 **Fig. 3. PIPs have significantly more mammalian adaptation than matched control proteins. (a)**

3 PIPs have evolved in great apes under similar levels of constraint as other proteins, as measured by

4 extant polymorphism. (b) PIPs have 39% more non-synonymous substitutions, relative to

5 synonymous substitutions, across 24 mammal species. (c) BUSTED detects recurrent adaptation in a

6 greater fraction of PIPs than in sets of matched controls. (d) BS-REL tests identify PIPs as evolving

7 adaptively on more branches in the mammalian tree. (e) BS-REL tests identify a higher proportion of

8 codons in PIPs as evolving adaptively. (f) The excess of adaptive codons in PIPs increases ($p=0.001$)

9 as the BUSTED threshold for including BS-REL estimates becomes more stringent (Methods VII). The

10 solid line indicates the mean excess; the dashed line indicates the 1:1 expectation; gray shading

11 indicates 95% confidence intervals. In all panels, *** = $p<0.001$; ns = $p>0.05$.

12

13

14 various metrics of function and genomic context (S4 Fig), we noted that PIPs have

15 the typical ratio of non-synonymous to synonymous polymorphism in great apes

16 (Fig 3A; mean $pN/(pS+1) = 0.21$ vs 0.20 ; $p=0.40$; see Enard et al., 2016). That is, PIPs

17 do not appear more or less evolutionarily constrained than other proteins,

1 bolstering the null expectation that they should evolve at average rates.

2 In contrast, we find that PIPs do have a significantly higher ratio of non-
3 synonymous to synonymous substitutions across 24 mammal species (Fig 3B; mean
4 $dN/dS = 0.195$ vs 0.140 , $p < 10^{-5}$). This ~40% elevation in dN/dS , despite
5 unremarkable $pN/(pS+1)$, conservatively suggests that over a third of substitutions
6 in PIPs may be related to positive selection imposed by malaria (or related
7 pressures).

8 We refined this result using the BS-REL and BUSTED tests (Methods VII),
9 available in the HYPHY package (Kosakovsky Pond et al., 2011; Murrell et al., 2015;
10 Pond et al., 2005). Both tests use maximum likelihood models to estimate the
11 proportion of codons in a protein with $dN/dS > 1$, consistent with adaptation. For
12 each protein, BUSTED estimates this adaptive proportion across the entire
13 phylogenetic tree, while BS-REL estimates it per branch

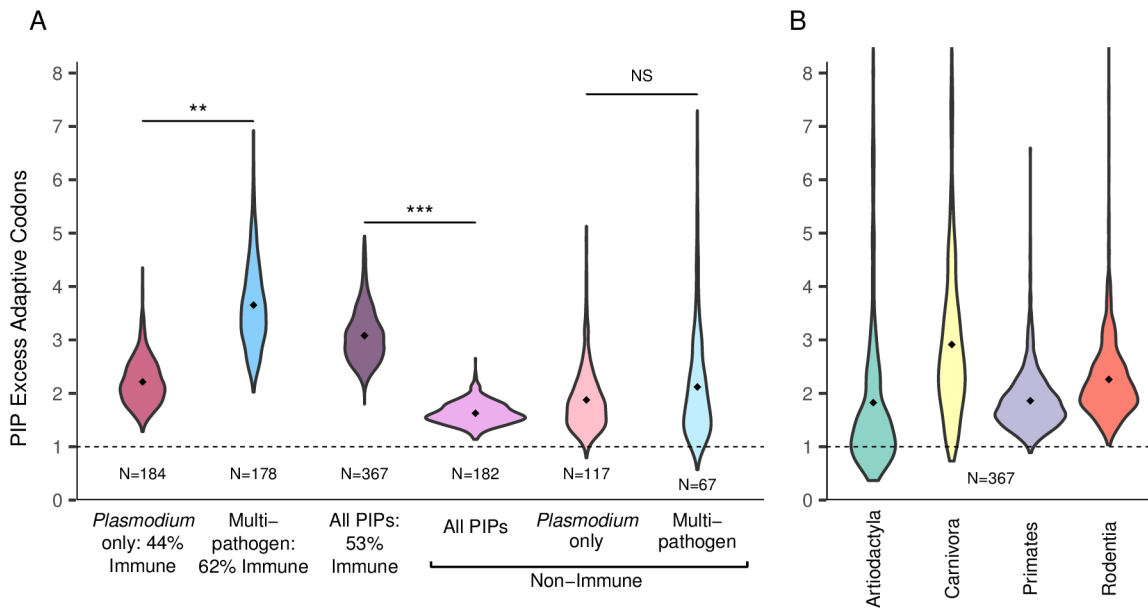
14 Both models find evidence of excess adaptation in PIPs. Over 37% of PIPs
15 have BUSTED evidence (at $p \leq 0.05$) of recurrent adaptation in mammals, versus 23%
16 of matched controls ($p < 10^{-5}$; Fig 3C). Similarly, PIPs have BS-REL evidence for
17 adaptation on more branches of the mammalian tree ($p = 1.87 \times 10^{-4}$; Fig 3D), and for
18 more codons per protein ($p < 10^{-5}$; Fig 3E). This excess is robust to the BUSTED p -
19 value threshold used to define adaptation, and increases as the threshold becomes
20 more stringent (Fig 3F, $p = 0.001$). Overall, these matched tests show that PIPs have
21 indeed experienced an accelerated rate of adaptive substitutions, consistent with
22 malaria as an important selective pressure.

23

1 ***High rate of adaptation in PIPs known to interact only with Plasmodium***

2 We have shown that a large set of host proteins with strong connections to
3 *Plasmodium* (STable 1, Fig 1 A-C) have, over deep time scales, evolved under
4 exceptionally strong positive selection (Fig 3). Given that nearly half of PIPs are
5 known to also interact with viruses and/or bacteria (Fig 1D), one critical question is
6 whether *Plasmodium* is truly the source of this selection. We attempted to isolate
7 *Plasmodium* as a selective pressure by dividing PIPs into '*Plasmodium*-only' and
8 'multi-pathogen' categories, based on the available information regarding viruses
9 and bacteria (Fig 1D; Methods IV). We find that *Plasmodium*-only PIPs have a 2.2-
10 fold excess of adaptation compared to matched controls ($p=0.008$; Fig 4A, far left),
11 when adaptation is measured as the proportion of adaptive codons per gene (Fig
12 3E). This suggests that *Plasmodium* may have specifically driven adaptation in a
13 large number of mammalian proteins, apart from any pleiotropic interactions they
14 may have with other pathogens.

15 Nonetheless, multi-pathogen PIPs have 3.7 times more adaptation than
16 matched controls—significantly higher than the excess in *Plasmodium*-only PIPs
17 ($p=0.005$; Fig 4A, left). This suggests that an increased number and diversity of
18 pathogen interactions may drive a cumulative increase in host adaptation.
19 Importantly, however, these multi-pathogen interactions are concentrated in
20 immune PIPs (Fig 1D; Fig 4A). Since immune genes are well known to evolve at
21 accelerated rates (Hurst and Smith, 1999; Nielsen et al., 2005; Bustamante et al.,
22 2005; Voight et al., 2006; Williamson et al. 2007; Sackton et al., 2007), this immune



1
2 **Fig. 4. PIP adaptation is pleiotropic and pervasive.** (a) Left: PIPs that interact with *Plasmodium*
3 alone have 2.2 times more adaptation than expected. PIPs that also interact with another viral and/or
4 bacterial pathogen have 3.7 times more adaptation than expected. Middle: Non-immune PIPs have
5 less adaptation than all PIPs. Right: The difference between Plasmodium-only and multi-pathogen
6 PIPs disappears when only non-immune PIPs are considered. All tested groups of PIPs exceed the 1:1
7 null expectation (dashed line) with $p < 0.05$. Each violin represents the ratio between PIPs and
8 matched controls over 1000 iterations, with the black point indicating the mean. Key: *** = $p < 0.001$;
9 ** = $p < 0.01$; NS=not significant. (b) Evidence of excess adaptation in PIPs is observed in all tested
10 orders of the mammalian tree. The ratio is significantly higher than 1 ($p < 0.05$) in all orders but
11 Artiodactyla ($p = 0.28$).
12
13
14 enrichment could confound the excess of adaptation observed in multi-pathogen
15 PIPs.

16 Before disentangling this issue, we first verified the correlation between
17 immune function and adaptation (Methods VI). We find that while PIPs overall have

1 adapted at a 3.1-fold higher rate than matched controls, non-immune PIPs have
2 adapted at a 1.7-fold higher rate than matched, non-immune controls (Fig. 4A,
3 middle). This difference, which is highly significant ($p < 0.001$), reinforces that
4 immune enrichment could confound adaptation in multi-pathogen PIPs. To isolate
5 these two effects, we then considered only non-immune PIPs, divided into groups by
6 their total number of pathogen interactions (Fig 4A, right; S7 Fig). In these non-
7 immune PIPs, in contrast to all PIPs, we find that additional interactions beyond
8 *Plasmodium* have no additional effects on adaptation.

9 Together, these results suggest that adaptation in immune genes is difficult to
10 attribute to any single selective pressure. The immune system appears to be the
11 most efficient avenue for hosts to simultaneously adapt to multiple pathogens. In
12 contrast, host adaptation to *Plasmodium* is apparent through both immune and non-
13 immune pathways (Fig 1D; Fig 4A). We have shown that non-immune genes evolve
14 more slowly and have less pathogen pleiotropy (Fig 4A; Fig 1D). Thus, though
15 *Plasmodium* has likely played a major role in immune evolution, we can be more
16 confident that selection imposed by *Plasmodium* has specially driven adaptation in
17 non-immune PIPs.

18

19 ***PIP adaptation is related to expression in blood, liver, and lung***

20 Malaria infections are biologically complex, and host adaptation to
21 *Plasmodium* could occur in genes expressed in several malaria-relevant tissues (Fig
22 1B). We used multiple linear regression to test whether the rate of adaptation in a
23 gene, as measured by BS-REL and BUSTED, was related to its tissue-specific

1 expression, as measured by GTEx

2 For PIPs, rates of adaptation are significantly and positively related to
 3 relative expression in blood, liver, and lung, but not in other malaria-related tissues
 4 (Table 1, column 2). Overall, in a multiple linear model, PIP expression in these
 5 tissues explains 17.4% of the variance in the proportion of adaptive codons. In
 6 contrast, the tissue-specific expression of matched control genes (Methods III)
 7 explains only 4.6% of this variance in adaptation, or 3.8 times less ($p < 0.001$). When
 8 compared to samples of control genes matched for total expression, as well as for
 9 expression in blood, liver, and lung, PIP relationships between adaptation and tissue
 10 expression are significantly stronger than expected (Table 1, column 3). This
 11 suggests that blood, liver, and lung, among all sites of PIP expression, may
 12 experience the strongest selective pressures from *Plasmodium* parasites.

13

14 **Table 1.** PIP adaptation is linked to tissue-specific expression.

Tissue	p-value, PIP adaptation ~ expression	p-value, PIP vs. matched controls
blood	7.90×10^{-3}	0.012
liver	9.47×10^{-16}	<0.001
lung	3.29×10^{-4}	<0.001
lymphocytes	0.378	0.488
spleen	0.193	0.419
s. adipose	0.658	0.584
v. adipose	0.380	0.302

15

16 Relative expression in the blood, liver, and lung is more strongly related to the proportion of adaptive
 17 codons in PIPs than in other genes. 1000 sets of matched controls were matched to PIPs for total
 18 expression, as well as for relative abundance of expression in blood, liver, and lung. Abbreviations: s –
 19 subcutaneous; v – visceral.

1 ***PIP adaptation is not limited to Plasmodium-infected lineages***

2 A number of *Plasmodium* species infect mammalian hosts in the orders
3 Primates and Rodentia (Carlton, Perkins, and Deitsch, 2013). In contrast,
4 Artiodactyla and Carnivora are parasitized by other genera of Apicomplexan
5 parasites, which also reproduce in the blood and are transmitted by insects (Clark
6 and Jacobson, 1998). To further test the specificity of PIP adaptation, we applied the
7 BUSTED and BS-REL models to separate protein alignments for each mammalian
8 order (Methods VIII).

9 When all PIPs are considered, we find significant excesses of adaptation in
10 rodents ($p < 0.001$), primates ($p = 0.005$), and carnivores ($p = 0.02$; Fig 4B). The signal
11 is positive, but not significant, in artiodactyls ($p = 0.28$; Fig 4B). Artiodactyls are the
12 most poorly-represented group in our mammalian tree (S1 Fig), and we observe
13 that noise in these estimates is negatively correlated with the number of sampled
14 species ($R^2 = 0.65$; S8 Fig). The signal also becomes noisier, especially for carnivores,
15 when the analysis is limited to *Plasmodium*-only PIPs (S9 Fig). Overall, we find no
16 statistical evidence for the restriction of PIP adaptation to certain mammal lineages,
17 consistent with the phenotypic similarity between *Plasmodium* and other
18 Apicomplexan pathogens. Other ubiquitous pathogens that interact with PIPs,
19 namely viruses and bacteria (Fig 1D), may further contribute to these mammal-wide
20 patterns.

21

22 ***Understanding a single case of adaptation to Plasmodium***

23 We have shown that *Plasmodium* has driven, at least in part, an accelerated

1 rate of adaptation in a set of 410 mammalian PIPs. In order to understand this
2 adaptation at a more mechanistic level, we selected a single PIP for more detailed
3 investigation.

4 Of the top ten PIPs with the strongest BUSTED evidence of adaptation, only
5 one candidate—alpha-spectrin, or *SPTA1*—has been extensively characterized for
6 molecular interactions with *Plasmodium* proteins. Alpha-spectrin is a textbook
7 example of a major structural component of the red blood cell (RBC) membrane. In
8 humans, several polymorphisms in this gene are known to cause deformations of
9 the RBC, which may either be symptomless or cause deleterious anemia (reviewed
10 in, e.g., Gallagher, 2004). The *SPTA1* protein has a well-defined domain structure,
11 and specific interactions with *Plasmodium* proteins are known for three domains
12 (Fig 5). Repeat 4 is the binding site for KAHRP, the major *P. falciparum* component
13 of the adhesive 'knobs' that form on the surface of infected RBCs (Pei et al., 2005).
14 Another 65-residue fragment containing EF-hand 2 has been shown to bind to
15 PfEMP3; this interaction destabilizes the RBC skeleton, potentially allowing mature
16 merozoites to egress from the cell (Pei et al., 2007). A central SH3 domain can also
17 be cleaved by a promiscuous *Plasmodium* protease called plasmepsin II (Le Bonniec
18 et al., 1999), which mainly functions in hemoglobin digestion (Francis et al., 1997).
19 Furthermore, naturally occurring mutations in the first three *SPTA1* domains have
20 been shown to impair the growth of *P. falciparum* in human RBCs (Schulman *et al.*
21 1990; Facer, 1995; Dhermy et al. 2007).

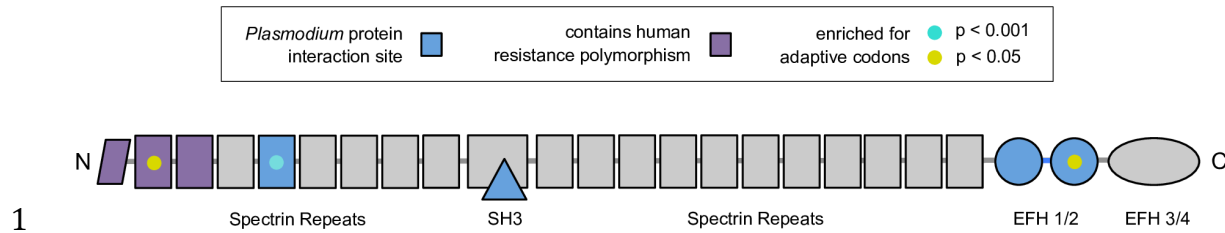


Fig 5. Domains of the alpha-spectrin (*SPTA1*) protein that are enriched for mammalian adaptation overlap with *Plasmodium* interaction sites. Adaptive codons were determined with MEME on an 85-species alignment of *SPTA1* coding sequence (Methods IX). *Plasmodium* interaction sites, and sites of human resistance mutations, were drawn from the literature (see text).

We wished to test whether sites of mammalian adaptation in *SPTA1* mapped to any of these *Plasmodium*-relevant domains. To identify adaptive codons with higher precision and power, we aligned *SPTA1* coding sequences from 61 additional mammal species (S5 Table) for analysis in MEME (Murrel et al., 2012; Methods IX). Of the 2,419 codons in this large mammalian alignment, we found that 63 show strong evidence of adaptation ($p < 0.01$), and that these are distributed non-randomly throughout the protein.

Remarkably, three domains—Repeat 1, Repeat 4, and EF-hand 2—are significantly enriched for adaptive codons, after controlling for domain length and conservation (Fig 5; Methods IX). That is, all three *SPTA1* domains with strong evidence of adaptation in mammals are known to either interact specifically with *P. falciparum* proteins, or harbor human mutations that provide resistance to *P. falciparum*. This overlap is unlikely to occur by chance ($p = 0.015$), and is robust to the p-value thresholds chosen (S6 Table). Thus, evidence from *SPTA1* suggests a

1 meaningful and specific connection between host adaptation and the mechanics of
2 *Plasmodium* infection.

3

4 Discussion

5 In this work, we have examined decades of malaria literature to expand the
6 collection of mammalian, *Plasmodium*-interacting proteins by over an order of
7 magnitude (Fig 1). We show that, compared to control proteins matched for various
8 properties (Fig 2), these 410 PIPs have adapted at exceptionally high rates in
9 mammals (Fig 3). The highest rates of adaptation are evident in immune PIPs,
10 especially those that share interactions with viruses and bacteria (Fig 4A). However,
11 we show that *Plasmodium* itself (or related Apicomplexans) has likely been an
12 important driver of this adaptation, especially for non-immune proteins (Fig 4A).

13 We used collections of available data on other pathogens to isolate a set of
14 PIPs that, to the best of our knowledge, lack any 'multi-pathogen' interactions. These
15 '*Plasmodium*-only' PIPs, whether immune or not, have adapted at over twice the
16 expected rate in mammals (Fig 4A). This suggests that *Plasmodium* has had an
17 appreciable effect on PIP evolution, beyond the effect of unrelated pathogens. Still,
18 many interactions with other pathogens likely remain unknown, making it
19 difficult—based on this evidence alone—to dismiss their importance.

20 However, two other pieces of evidence support *Plasmodium* as a key selective
21 pressure. First, mammal-wide adaptation in PIPs is strongly linked to PIP
22 expression in human blood, liver, and lung (Table 1). *Plasmodium* parasites are well
23 known to replicate within red blood cells (RBCs) and hepatocytes, and infected

1 RBCs tend to sequester in the lungs, with serious consequences (e.g. Aursudkij et al.,
2 1998). Thus, the pathophysiology of malaria is reflected in the tissues where PIPs
3 show the strongest evidence of adaptation.

4 Second, in the well-studied case of alpha-spectrin, we show that domain-level
5 interactions with *Plasmodium* perfectly explain the observed patterns of adaptive
6 substitution (Fig 5). Besides validating the ability of codon evolution models to
7 detect adaptation at particular residues (Methods VII), this result affirms a specific
8 role for *Plasmodium* in mammalian evolution, beyond the immune-focused role
9 played by pathogens in general (Fig 4A). Thus, despite the inevitability of at least
10 some pleiotropy (Wagner and Zhang, 2011), we show that phenotypic information
11 can be leveraged to link genetic adaptation to specific sources of selection.

12 Throughout this work, we showcase the utility of phenotypic information for
13 studying evolution. We demonstrate that recent, well-funded projects like GTEx and
14 Encode can provide, among many other uses, the raw information required for
15 meaningful evolutionary comparisons (Fig 2; Methods). Smaller-scale projects,
16 including most of the scientific papers contained in PubMed, also contain an
17 impressive quantity of valuable data (Fig 1A). However, we find that heavy manual
18 curation is still required to remove false positives from literature searches (Methods
19 I). In the future, unique and automatic indexing of existing data will be key to
20 understanding the evolution of complex phenotypes, and should be a major
21 research focus, alongside the accessibility of new data.

22 Finally, this work provides an interesting contrast with previous studies,
23 which have associated only a few dozen human genes with malaria resistance

1 (Verra et al., 2009). Only a handful of these genes are backed by convincing evidence
2 of positive selection in humans, and nearly all of these are RBC proteins (Hedrick,
3 2011; MalariaGEN, 2015). In contrast, our work provides a repository of hundreds
4 of diverse human genes with phenotypic links to malaria (Fig 1; S1 Table). Why,
5 then, do we know of so few examples of recent human adaptation to *Plasmodium*?
6 This disconnect may depend simply on the timescale of human evolution, which is
7 only a fraction of the ~105 million years of mammalian evolution (Murphy et al.,
8 2007). Or, perhaps the difficulty of detecting balancing selection (Charlesworth,
9 2006) has obscured additional, important human variants. Future work will utilize
10 the large set of PIPs to better understand the evolution of malaria resistance in
11 humans.

12 In conclusion, we have found evidence of substantially accelerated
13 adaptation in mammalian proteins that interact with *Plasmodium*. In the case of
14 rapidly evolving immune proteins, *Plasmodium* appears to share responsibility with
15 other groups of pathogens, including viruses and bacteria. We show that it can be
16 difficult to attribute evolutionary changes to a single selective agent, given the
17 surprising pleiotropy of host genes with regard to very different pathogenic agents.
18 But in many cases—as well as in the case of alpha-spectrin—our approach allows us
19 to infer that *Plasmodium*-like parasites have imposed a substantial selective
20 pressure on mammals. We hope that our collection of 410 mammalian PIPs will
21 continue to prove a powerful resource for exploring host interactions with
22 *Plasmodium*.

1 **Methods**

2 *I. Identification of PIPs*

3 We queried PubMed for scientific papers containing both a gene name and
4 the word(s) 'malaria,' '*falciparum*,' or '*Plasmodium*' in the title or abstract, as of May
5 21, 2015. Human gene names were drawn from the HUGO Gene Nomenclature
6 Committee (Gray et al., 2015; <http://www.genenames.org/>) for 9,338 mammalian
7 orthologs (Methods II). For each of the 2,249 genes that returned a hit, we manually
8 evaluated the titles of up to 20 associated papers to assess the link between the gene
9 and a malaria phenotype. Many acronyms used to represent genes are also used as
10 abbreviations for techniques, locations, drugs, or other phrases. Consequently, most
11 genes could be eliminated based on their nominal connection with papers
12 addressing non-genetic aspects of malaria.

13 For papers discussing genes, we examined the abstracts for the presence and
14 type of evidence connecting genes to malaria phenotypes. In cases where the
15 abstract was ambiguous, we examined the full text of the paper. To limit the number
16 of false positives, we did not include results from RNAseq or other high-throughput
17 experiments.

18

19 *II. Generation of mammalian ortholog alignments*

20 We used BLAT to identify homologs of 22,074 human coding sequences in 24
21 high-depth mammal genomes (S1 Fig). We retained orthologs which (1) had best
22 reciprocal hits in all 24 mammal species, (2) lacked any in-frame stop codons, (3)
23 were at least 30% of the length of the human sequence, and (4) had clearly

1 conserved synteny in at least 18 non-human species. Coding sequences for the
2 resulting 9,338 proteins were aligned with PRANK, and any codon present in fewer
3 than eight species was excluded from analysis. Additional details are available in
4 Enard et al. (2016).

6 *III. Tissue Expression Analyses*

7 Expression data for 53 human tissues were downloaded from the GTEx
8 portal (<http://www.gtexportal.org/home/>) on October 18, 2015. For tissue-specific
9 analyses (Fig 1B, Table 1), we converted RPKM values to relative abundance (RA)
10 values for each tissue. RA is simply the proportion of each gene's total RPKM found
11 in each tissue. For matching controls, we summed RPKM values over all tissues to
12 yield total expression.

13 Because PIPs have substantially higher total expression than other proteins
14 (Fig. 2A), for each PIP, we identified matched control proteins with +/-20% of the
15 total RPKM expression of that PIP. For each of 10,000 iterations, we randomly
16 selected one matched control out of the set of potential matches for each PIP,
17 producing PIP and control sets of equal size and with indistinguishable distributions
18 of total expression (S4 Fig). We compared the mean value from each matched
19 control set to the mean value for PIPs, and determined p-values empirically, as the
20 fraction of permutations with control mean \geq PIP mean. Tissues with significantly
21 different PIP expression were determined after applying the Bonferroni correction
22 for 53 tissue tests (S2 Fig).

1 We also correlated tissue-specific expression in malaria tissues with the
2 proportion of adaptive codons per gene (Table 1), using multiple linear regression
3 in R. To generate matched controls for this analysis, we matched control genes to
4 PIPs based on their RA in blood, liver, and lung, as well as total expression.

5

6 *IV. Collection of VIPs and BIPs*

7 Virus-interacting proteins (VIPs) were manually curated by Enard *et al.*
8 (2016), in the same manner as PIPs. To our knowledge, no similar collection of high-
9 quality interactions is available for other pathogens. Therefore, we queried the EBI
10 IntAct database (<http://www.ebi.ac.uk/intact/>) for protein interactions between
11 Kingdom Bacteria (taxid:2) or Phylum Apicomplexa (taxid:5794) and humans
12 (taxid:9606). This approach, while much faster than manual curation, is less ideal
13 for two reasons: (1) many interactions are not included in the database (e.g., only 17
14 human-*Plasmodium* interactions are included in IntAct), and (2) many of the
15 included interactions are based on high-throughput assays, including yeast two-
16 hybrid experiments, which suffer from both false negatives and false positives
17 (Brückner et al 2009). Consequently, we do not perform rigorous analysis for
18 bacterial-interacting proteins (BIPs), as has been done for PIPs and VIPs (Enard et
19 al., 2016). Rather, we use the IntAct BIPs only to classify PIPs as '*Plasmodium*-only'
20 or 'multi-pathogen.'

21

22 *V. Collection of other protein metrics*

23 GO annotations were downloaded in October, 2015 from the Gene Ontology

1 website (Ashburner et al., 2000; <http://geneontology.org/>)

2 Regions of DNaseI hypersensitivity, combined from 95 cell types, were
3 obtained from the databases of the ENCODE Project Consortium (2012;
4 <https://www.encodeproject.org/>). We calculated the density of DNaseI
5 hypersensitivity regions in 50 Kb windows centered on each ortholog.

6 Protein expression levels were obtained from the Human Proteome Map
7 (Kim et al., 2014; <http://www.humanproteomemap.org/>), which used high
8 resolution and high accuracy Fourier transform mass spectrometry experiments.
9 We summed spectral values over 30 tissues and cell types and took the log of these
10 total values. The log number of interacting partners for each human protein was
11 obtained from the Biogrid Database (Stark et al., 2011; <http://thebiogrid.org/>),
12 curated by Luisi et al., 2015.

13 Genomic elements conserved in 46 vertebrate species, derived from
14 PhastCons (Siepel et al., 2005), were downloaded from the UCSC genome browser
15 (<http://hgdownload.cse.ucsc.edu/goldenPath/hg19/phastCons46way/>). We
16 calculated conserved element density within 50 kb windows centered on each gene
17 in the human reference. Coding density was calculated from coding nucleotides in
18 the same 50 Kb windows. The length and GC content of each protein was derived
19 from the mammalian alignment (Methods II).

20 We assessed citation frequency of each gene in two ways. First, we counted
21 the citations linked to each gene on its PubMed Gene page
22 (<http://www.ncbi.nlm.nih.gov/gene>). Second, we downloaded the Gene References
23 into Function, or GeneRIFs, contributed to PubMed by scientists

1 (<ftp://ftp.ncbi.nih.gov/gene/GeneRIF/>). These measures were highly correlated
2 (not shown), and only citations are reported.

3

4 *VI. Matching PIPs to control proteins*

5 Each PIP was matched to a set of control proteins based on similarity in five
6 metrics: mRNA expression, protein expression, protein-protein interactions, DNaseI
7 density, and conserved element density (Fig. 2A-D; G). We allowed a control protein
8 to be considered a PIP 'match' if each of its five values fell within a given range,
9 based on the PIP values. For example, margins of min=0.1 and max=0.2 for mRNA
10 expression means that, for a control protein to be matched to a PIP, the mRNA
11 expression of the control must fall between 90-120% of the mRNA expression of the
12 PIP. We wished to maximize the number of matched controls per PIP while creating
13 control sets that were statistically indistinguishable from PIPs for all five metrics
14 (e.g. S4 Fig). To achieve this balance, we iteratively chose the maximum margins
15 that yielded average p-values, over 100 permutations, of ≥ 0.1 for each metric. Once
16 appropriate margins were found, we obtained matched control sets of equal size to
17 the PIP set by randomly sampling one matched control protein for each PIP. For
18 each permutation test, 10,000 sets of matched controls were sampled.

19 Margins for the main permutation test (Fig 3) are given in S4 Table. For
20 subsets of PIPs (i.e., Fig 4A, S7 Fig), the margins were altered to generate well-
21 matched controls in every case. For Table 1, only the stated expression values and
22 $pN/(pS+1)$ were checked for matching, to avoid excluding too many PIPs. Because
23 we chose sets of matched controls based on the distribution of PIP values included

1 in each test, whether any given PIP was matched depended on the other PIPs in the
2 test (e.g., one extreme PIP may or may not be balanced out by another). Therefore,
3 the sum of matched PIPs across categories differs slightly from the total (Fig 4A).

4 The pool of immune controls is relatively small (998 genes), compared to the
5 pool of non-immune controls with other GO annotations (7,594 genes)(S2 Table).

6 This made it difficult to match immune PIPs to immune controls, without discarding
7 many immune PIPs. Consequently, to test hypotheses of accelerated immune
8 adaptation, we focused on comparing all PIPs to all controls and non-immune PIPs
9 to non-immune controls (Fig 4A). For these and other tests in Figure 4, 1000 sets of
10 matched controls were sampled for each violin.

11

12 *VII. Estimating adaptation with models of codon evolution*

13 We used the codeml model m8 from the PAML package (Yang, 2007) to
14 estimate dN/dS for each ortholog (Fig 3B). However, branch-site tests in PAML rely
15 on assumptions that may be violated in the case of recurrent adaptation to a
16 pervasive selective pressure (see Enard et al. 2016). Consequently, we chose to
17 implement the maximum-likelihood branch-site tests in the better-performing
18 HYPHY package (Kosakovsky Pond et al., 2011). We used the BUSTED algorithm
19 (Murrell et al., 2015) to detect recurrent selection across the entire tree for each
20 gene, and BS-REL to estimate the proportion of positively selected codons in each
21 gene on each branch. Both of these algorithms rely on the same underlying codon
22 model; details of the model are described in Kosakovsky Pond et al. (2011), Murrell
23 et al. (2015), and reviewed in Enard et al. (2016). Unless otherwise specified (i.e.,

1 Fig 3F), codons identified by BS-REL were 'counted' as adaptive if the BUSTED p-
2 value for that gene was ≤ 0.05 . For each evolutionary statistic (i.e. adaptive codons,
3 adaptive branches, dN/dS, pN/pS), empirical p-values were derived by comparing
4 the average value of the PIPs to the average values of 100,000 permutations of
5 matched controls.

6

7 *VIII. Order-specific analyses*

8 We split the mammal-wide alignments for each gene into four non-
9 overlapping alignments, corresponding to the following clades: **primates** (human,
10 chimp, gorilla, orangutan, gibbon, macaque, baboon, marmoset, bushbaby), **rodents**
11 (mouse, rat, guinea pig, squirrel, rabbit), **carnivores** (panda, ferret, dog, cat), and
12 **artiodactyls** (sheep, cow, pig) (see S1). We excluded microbat, elephant, and horse,
13 as these species are not closely related to any of the four major groups (Murphy et
14 al., 2007; S1 Fig). However, we included rabbit with rodents, because they are more
15 closely related. We ran BUSTED on each alignment to yield a p-value of clade-
16 specific adaptation for each gene.

17 PIPs were matched to controls as described above (Methods V). However,
18 rather than counting BS-REL adaptive codons in all branches if the tree-wide
19 BUSTED $p \leq 0.05$, we (1) kept each clade codon count separate, (2) counted codons
20 only on branches within a clade, and (3) counted codons only if the clade-specific
21 BUSTED $p \leq 0.05$. The ratio of adaptive codons in PIPs versus controls was then
22 calculated as before, by taking 1,000 random samples of matched controls.

23

1 IX. Alpha-spectrin

2 Alpha-spectrin homologs were initially identified in 88 mammal species
3 using NCBI Gene (http://www.ncbi.nlm.nih.gov/gene/?Term=ortholog_gene_6708).
4 The sequence of the longest mRNA transcript for each species was downloaded
5 using E-Utilities, and each transcript was trimmed to the longest ORF using
6 TransDecoder (Haas et al., 2013; <http://transdecoder.github.io/>). Coding sequences
7 with <50% of the human CDS length were removed. The remaining 85 coding
8 sequences were aligned with PRANK (Löytynoja and Goldman, 2008) using default
9 settings (S5 Table). The alignment was manually inspected and corrected using
10 JalView (Waterhouse et al., 2009).

11 A phylogenetic tree for the 85 species was obtained, using NCBI Taxonomy,
12 from phyloT (<http://phylot.biobyte.de/>). This tree, along with the corrected
13 alignment, was inputted into HyPhy to run MEME (Murrell et al., 2012), which
14 yielded a p-value of adaptation for each codon. We used the domain designations
15 from SMART (Schultz et al., 1998; <http://smart.embl-heidelberg.de/>) to assign
16 92.2% of *SPTA1* codons to one of 25 domains (S6 Table). Then, for each domain, we
17 calculated an 'adaptation score' as:

$$18 \quad a/v$$

19 where a measures adaptation (the proportion of codons within the domain with
20 MEME $p \leq 0.01^*$) and v measures variability (the proportion of codons within the
21 domain that vary among species, i.e., are not 100% conserved). This score also
22 controls for domain length, as it uses the proportion of codons within the domain.
23 To calculate the significance of each domain's adaptation score (i.e., to ask, is it

1 higher than expected?), we randomly permuted codons among domains 10,000
2 times.

3 *We also tested MEME p-value cutoffs of 0.1, 0.5, 0.005, and 0.001 for
4 defining α ; these results are available in S6 Table. The results for $p \leq 0.01$, which are
5 reported in the main text, are representative across these cutoffs.

6

7 **Data Access**

8 All data used in this work are publicly available (Methods I-V). The collection of PIPs
9 is available in S1 Table.

10

11 **Acknowledgements**

12 We wish to thank Kerry Geiler-Samerotte for her thoughtful comments on the
13 manuscript, along with the rest of the Petrov lab. ERE thanks Jane Carlton for
14 abbreviation advice; Jamie Blundell and Anisa Noorassa for figure advice; and Daniel
15 Friedman, for conceding that 'protein' can mean 'gene.'

16

17 **Author Contributions**

18 E.R.E curated the PIPs. D.E., E.R.E., and N.T. collected other data. E.R.E. and D.E.
19 performed the analyses, with design input from D.A.P. and S.V. E.R.E. and D.A.P.
20 wrote the paper, with contributions from all other authors.

21 This work was supported by NIH grants R01GM089926 and R01GM097415 and
22 NSF grant R35GM118165-01 to DAP, and an NSF Graduate Research Fellowship to

1 ERE (DGE-1247312).

2

3 **References**

4 Aidoo M, Terlouw DJ, Kolczak MS, McElroy PD, ter Kuile FO, Kariuki S, et al.

5 Protective effects of the sickle cell gene against malaria morbidity and

6 mortality. *Lancet*. 2002;359: 1311–1312.

7 Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, et al. Gene ontology:

8 tool for the unification of biology. *Nat Genet*. 2000;25(1): 25-9.

9 Aursudkij B, Wilairatana P, Vannaphan S, Walsh DS, Gordeux VR, Looareesuwan S.

10 Pulmonary edema in cerebral malaria patients in Thailand. *Southeast Asian J*

11 *Trop Med Public Health*. 1998;29(3): 541-5.

12 Brückner A, Polge C, Lentze N, Auerbach D, Schlattner U. Yeast Two-Hybrid, a

13 Powerful Tool for Systems Biology. *Int J Mol Sci*. 2009; 10(6): 2763-88.

14 Bustamante CD, Fledel-Alon A, Williamson S, Nielsen R, Hubisz MT, Glanowski S,

15 et al. Natural selection on protein-coding genes in the human genome.

16 *Nature*. 2005;437: 1153-1157.

17 Carlton JM, Perkins DJ, Deitsch KW, editors. *Malaria Parasites: Comparative*

18 *Genomics, Evolution and Molecular Biology*. 1st ed. Norfolk, UK: Caister

19 Academic Press; 2013.

20 Charlesworth D. Balancing Selection and Its Effects on Sequences in Nearby

21 Genome Regions. *PLOS Genetics*. 2006;2(4): e64. doi:

22 10.1371/journal.pgen.0020064

- 1 Clark IA, Jacobson LS. Do babesiosis and malaria share a common disease
2 process? *Annals of Tropical Medicine and Parasitology*. 1998; 92(4): 483-8.
- 3 Coatney GR, Roudabush RL. A Catalog and Host-Index of the Genus *Plasmodium*.
4 *J Parasitol*. 1936;22(4): 338-353.
- 5 Daub JT, Hofer T, Cutivet E, Dupanloup I, Quintana-Murci L, Robinson-Rechavi M,
6 Excoffier L. Evidence for polygenic adaptation to pathogens in the human
7 genome. *Mol Biol Evol*. 2013;30(7): 1544-58.
- 8 Dhermy D, Schrével J, Lecomte M-C. Spectrin-based skeleton in red blood cells and
9 malaria. *Curr Opin Hematol*. 2007;14(3): 198-202.
- 10 Enard D, Cai L, Petrov DA. Viruses are a dominant driver of protein adaptation in
11 mammals. *eLIFE*. 2016;5: e12469. doi: <http://dx.doi.org/10.7554/eLife.12469>
- 12 ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the
13 human genome. *Nature*. 2012;489(7414): 57-74.
- 14 Engwerda CR, Beattie L, Amante FH. The importance of the spleen in malaria. *Trends in*
15 *Parasitology*. 2005;21(2): 75-80.
- 16 Escalante AA, Ayala FJ. Evolutionary origin of *Plasmodium* and other
17 Apicomplexa based on rRNA genes. *PNAS*. 1995;92(13): 5793-5797.
- 18 Facer CA. Erythrocytes carrying mutations in spectrin and protein 4.1 show
19 differing sensitivities to invasion by *Plasmodium falciparum*. *Parasitol Res*
20 1995;81: 52-57.
- 21 Francis SE, Sullivan DJ Jr, Goldberg DE. Hemoglobin metabolism in the malaria
22 parasite *Plasmodium falciparum*. *Annu Rev Microbiol*. 1997;51: 97-123.
- 23 Franke-Fayard B, Janse CJ, Cunha-Rodrigues M, Ramesar J, Büscher P, Que I, et al.

- 1 Murine malaria parasite sequestration: CD36 is the major receptor, but
2 cerebral pathology is unlinked to sequestration. PNAS. 2005;102(32): 11468
3 73.
- 4 Gallagher PG. Hereditary Elliptocytosis: Spectrin and Protein 4.1R. Seminars
5 in Hematology 2004;41(2): 142-164.
- 6 Gray K, Yates B, Seal RL, Wright MW, Bruford EA. Genenames.org: the HGNC
7 resources in 2015. Nucleic Acids Res. 2015; 43(Database issue): D1079
8 D1085.
- 9 GTEx Consortium. The Genotype-Tissue Expression (GTEx) pilot analysis:
10 Multitissue gene regulation in humans. Science. 2015;348(6235): 648-60.
- 11 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo
12 transcript sequence reconstruction from RNA Seq: reference generation and
13 analysis with Trinity. Nat Protol. 2013;8(8): 1494-512. doi:
14 10.1038/nprot.2013.084.
- 15 Hedrick PW. Population genetics of malaria resistance in humans.
16 Heredity. 2011;107(4): 283-304.
- 17 Hulsen T, de Vlieg J, Alkema W. BioVenn—a web application for the comparison and
18 visualization of biological lists using area-proportional Venn diagrams. BMC
19 Genomics. 2008;9: 488.
- 20 Hurst LD, Smith NG. Do essential genes evolve slowly? Current Biology. 1999;9(14):
21 747-50.
- 22 Idro R, Marsh K, John CC, Newton CR. Cerebral malaria: mechanisms of brain injury

- 1 and strategies for improved neurocognitive outcome. *Pediatric Research*.
2 2010; 68(4): 267-74
- 3 Kim MS, Pinto SM, Getnet D, Nirujogi RS, Manda SS, Chaerkady R, et al. A draft map
4 of the human proteome. *Nature*. 2014;509(7502): 575-81.
- 5 Kosakovsky Pond SL, Murrell B, Fourment M, Frost SDW, Delpont W, Scheffler K. A
6 Random Effects Branch-Site Model for Detecting Episodic Diversifying
7 Selection. *Mol Bio Evol*. 2011;28(11): 3033-3043.
- 8 Kwiatkowski DP. How Malaria Has Affected the Human Genome and What
9 Human Genetics Can Teach Us about Malaria. *Am J Hum Genet*. 2005;77(2):
10 171-192.
- 11 Le Bonniec S, Deregnacourt C, Redeker V, Banerjee R, Grellier P, Goldberg DE,
12 Schrével J. Plasmeprin II, an acidic hemoglobinase from the *Plasmodium*
13 *falciparum* food vacuole, is active at neutral pH on the host erythrocyte
14 membrane skeleton. *J Biol Chem*. 1999;274(20): 14218-23.
- 15 Lovegrove FE, Gharib SA, Peña-Castillo L, Patel SN, Ruzinski JT, Hughes TR, Liles WC,
16 Kain KC. Parasite Burden and CD36-Mediated Sequestration Are
17 Determinants of Acute Lung Injury in an Experimental Malaria Model. *PLoS*
18 *Pathogens*. 2008;4(5): e1000068. doi: 10.1371/journal.ppat.1000068
- 19 Löytynoja A, Goldman N. A model of evolution and structure for multiple
20 sequence alignment. *Philosophical Transactions of the Royal Society B:*
21 *Biological Sciences*. 2008;363: 3913-3919.

- 1 Luisi P, Alvarez-Ponce D, Pybus M, Fares MA, Bertranpetit J, Laayouni H. Recent
2 positive selection has acted on genes encoding proteins with more
3 interactions within the whole human interactome. *Genome Biology and*
4 *Evolution* 2015;7: 1141-1154.
- 5 MalariaGEN, Band G, Rockett KA, Spencer CC, Kwiatkowski DP. A novel locus of
6 resistance to severe malaria in a region of ancient balancing selection.
7 *Nature*. 2015;526(7572): 253-7.
- 8 Murphy WJ, Pringle TH, Crider TA, Springer MS, Miller W. Using genomic data to
9 unravel the root of the placental mammal phylogeny. *Genome Res*.
10 2007;17(4): 413-421.
- 11 Murrell B, Wertheim JO, Moola S, Weighill T, Scheffler K, Kosakovsky Pond SL.
12 Detecting Individual Sites Subject to Episodic Diversifying Selection. *PLOS*
13 *Genetics*. 2012;8(7): e1002764. doi: 10.1371/journal.pgen.1002764.
- 14 Murrell B, Weaver S, Smith MD, Wertheim JO, Murrell S, Aylward A, et al. Gene-Wide
15 Identification of Episodic Selection. *Mol Biol Evol*. 2015;32(5): 1365-71.
- 16 Nielsen R, Bustamante C, Clark AG, Glanowski S, Sackton TB, Hubisz MJ, et al. A
17 Scan for Positively Selected Genes in the Genomes of Humans and
18 Chimpanzees. *PLOS Biology*. 2005;3(6): e170.
- 19 Orchard S, Ammari M, Aranda B, Breuza L, Briganti L, Broackes-Carter F, et al. The
20 MIntAct project—IntAct as a common curation platform for 11 molecular
21 interaction databases. *Nucl Acid Res*. 2014;42(Database issue): D358-63.
- 22 Ozinsky A, Underhill DM, Fontenot JD, Hajjar AM, Smith KD, Wilson CB, Schroeder L,
23 Aderem A. The repertoire for pattern recognition of pathogens by the innate

- 1 immune system is defined by cooperation between Toll-like receptors. PNAS.
2 2000;97(25): 13766-71.
- 3 Pei X, An X, Guo X, Tarnawski M, Coppel R, Mohandas N. Structural and functional
4 studies of interaction between *Plasmodium falciparum* knob-associated
5 histidine-rich protein (KAHRP) and erythrocyte spectrin. J. Biol. Chem.
6 2005;280(35): 31166-71.
- 7 Pei X, Guo X, Coppel R, Mohandas N, An X. *Plasmodium falciparum* Erythrocyte
8 Membrane Protein 3 (PfEMP3) destabilizes erythrocyte membrane skeleton.
9 J Biol. Chem. 2007;282(37): 26754-8.
- 10 Piel FB, Patil AP, Howes RE, Nyangiri OA, Gething PW, Williams TN, Weatherall DJ,
11 Hay SI. Global distribution of the sickle cell gene and geographical
12 confirmation of the malaria hypothesis. Nat Commun. 2010;1(104). doi:
13 10.1038/ncomms1104
- 14 Pond SL, Frost SD, Muse SV. Hyphy: Hypothesis testing using phylogenies.
15 Bioinformatics. 2005;21: 676-679.
- 16 Sackton TB, Lazzaro BP, Schlenke TA, Evans JD, Hultmark D, Clark AG. Dynamic
17 evolution of the innate immune system in *Drosophila*. Nat Genet.
18 2007;39(12): 1461-8.
- 19 Schulman S, Roth EF Jr, Cheng B, Rybicki AC, Sussman II, Wong M, Wang W, Ranney
20 HM, Nagel RL, Schwartz RS. Growth of *Plasmodium falciparum* in human
21 erythrocytes containing abnormal membrane proteins. PNAS. 1990;87(18):
22 7339-7343.
- 23 Schultz J, Milpetz F, Bork P, Ponting CP. SMART, a simple modular architecture

1 research tool: identification of signaling domains. PNAS. 1998;95(11): 5857
2 64.

3 Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al.
4 Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
5 genomes. Genome Research. 2005;15(8): 1034-1050.

6 Stark C, Breitkreutz BJ, Chatr-Aryamontri A, Boucher L, Oughtred R, Livstone MS, et
7 al. The BioGRID Interaction Database: 2011 update. Nucleic Acids Research.
8 2011;39(Database issue): D698-D704.

9 Travisano M, Shaw RG. Lost in the Map. Evolution. 2012;67(2): 305-314.

10 Verra F, Mangano MD, Modiano D. Genetics of susceptibility to *Plasmodium*
11 *falciparum*: from classic malaria resistance genes towards genome-wide
12 association studies. Parasite Immunol. 2009;31(5): 234-53.

13 Voight BF, Kudaravalli S, Wen X, Pritchard JK. A Map of Recent Positive Selection in
14 the Human Genome. PLOS Biol. 2006;4(3): e72.

15 Wagner GP, Zhang J. The pleiotropic structure of the genotype–phenotype
16 map: the evolvability of complex organisms. Nat Rev Gen. 2011;12(3): 204
17 13.

18 Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ (2009) Jalview Version
19 2-a multiple sequence alignment editor and analysis workbench.
20 Bioinformatics. 2009;25: 1189-1191. doi:10.1093/bioinformatics/btp033

21 Welch SG, McGregor IA, Williams K. The Duffy blood group and malaria prevalence
22 in Gambian West Africans. Trans Royal Society of Trop Med and Hygiene.
23 1977;71(4): 295-6.

- 1 WHO (2015). World Malaria Report 2015. Available from:
- 2 <http://www.who.int/malaria/publications/world-malaria-report-2015/en/>
- 3 Williamson SH, Hubisz MJ, Clark AG, Payseur BA, Bustamante CD, Nielsen R.
- 4 Localizing Recent Adaptive Evolution in the Human Genome. PLOS Genetics.
- 5 2007;3(6): e90.
- 6 Yang Z. PAML 4: Phylogenetic Analysis by Maximum Likelihood. Mol Biol Evol.
- 7 2007;24(8):1586-91.
- 8 Yamamoto H, Fara AF, Dasgupta P, Kemper C. CD46: The 'multitasker' of
- 9 complement proteins. Int J Biochem Cell Biol. 2013;45(12): 2808-20.