

An antigenic diversification threshold for falciparum malaria transmission at high endemicity

Qixin He¹, Mercedes Pascual^{1,2*},

1 Department of Ecology and Evolution, University of Chicago, Chicago, IL 60637, USA

2 Santa Fe Institute, Santa Fe, NM, 87501, USA

* pascualmm@uchicago.edu

Abstract

In malaria and several other important infectious diseases, high prevalence occurs concomitantly with incomplete immunity. This apparent paradox poses major challenges to malaria elimination in highly endemic regions, where asymptomatic *Plasmodium falciparum* infections are present across all age classes creating a large reservoir that maintains transmission. This reservoir is in turn enabled by extreme antigenic diversity of the parasite and turnover of new variants. We present here the concept of a threshold in local pathogen diversification that defines a sharp transition in transmission intensity below which new antigen-encoding genes generated by either recombination or migration cannot establish. Transmission still occurs below this threshold, but diversity of these genes can neither accumulate nor recover from interventions that further reduce it. An analytical expectation for this threshold is derived and compared to numerical results from a stochastic individual-based model of malaria transmission that incorporates the major antigen-encoding multigene family known as *var*. This threshold we call R_{div} ; it is complementary to the one defined by the classic basic reproductive number of infectious diseases, R_0 , which does not easily apply under large and dynamic strain diversity. This new threshold concept can be exploited for effective malaria control and applied more broadly to other pathogens with large multilocus antigenic diversity.

Introduction

The reproductive number of an infectious disease, R_0 , quantifies its epidemic growth potential and provides a threshold condition for the spread and control of pathogens [1]. This number has been applied extensively to pathogens with either no antigenic variation or a typically low number of genetically defined strains that are relatively stable in space and time. It becomes problematic how to evaluate and even to define it when antigenic variation is large and dynamic, challenging the actual definition of a strain. This is particularly the case for *Plasmodium falciparum* as well as for other pathogens with multigene and multilocus encoding of antigens, combined with extensive recombination [2–5]. Furthermore, R_0 does neither predict the turnover rate of new genes encoding antigens, nor explain how transmission characteristics and accumulation of these new genes influence each other and responses to control.

Under high transmission of *P. falciparum*, the major antigen-encoding gene family of the blood stage of malaria infection, known as *var*, acquires vast antigenic diversity via large gene copy numbers as well as ectopic recombination. Laboratory experiments have shown that a naïve infection can generate about sixty new recombinants per year [6, 7]

although this has not yet been demonstrated to occur in nature. A parasite typically harbors 40 to 60 gene copies across its chromosomes, with a large pool of gene variants within local populations reaching the tens of thousands [8,9]. Thus, parasites share locally only a few common *var* genes between different strains [9–11], and across seasons [12]. Spatial diversity in *var* genes has also been documented [8,13,14] indicating that migration from surrounding areas also contributes to new diversity and to the immunological challenge.

The accumulation and turnover of new antigenic variants constitutes a major impediment to control in regions of high endemism, where it underlies the large reservoir of chronic asymptomatic infections that sustains transmission. Today the global burden of *Plasmodium falciparum* is concentrated in these high transmission endemic areas within fifteen countries, mainly in sub-Saharan Africa (WHO 2017). A similar reservoir is found in other vector-borne diseases that exhibit a high prevalence of infection with no clinical symptoms in domestic and wildlife hosts [15–17]. Nonsterile specific immunity is common to all these pathogens as a result of extreme antigenic variation encoded by multigene families [18,19].

We present here a reproductive number complementary to R_0 that defines a threshold for parasite antigenic diversification, below which the accumulation of new antigen-encoding genes no longer occurs even though they are consistently produced. We introduce the concept for infectious agents in general, derive an analytical expectation for the rate of generation of “successful” new genes for the *var* system in *P. falciparum*, and demonstrate the existence of the predicted analytical threshold in numerical simulations of a stochastic agent-based model that incorporates *var* genes and the acquisition of immunity by individual hosts. We then investigate the epidemiological and evolutionary factors that influence this diversification rate analytically. We show that this rate for the accumulation of genetic novelty, we call R_{div} , maps onto transmission intensity, separating at a threshold a regime in which new genes are able to accumulate from one in which they are unable to do so, despite transmission still occurring (i.e., R_0 remaining above one). We discuss implications for malaria control and elimination, future directions to estimate and monitor distance to this quantity in high transmission endemic regions, and its applicability to other infectious diseases.

Results

Consider genes of a parasite whose mutation generates new antigenic variation during transmission and infection. Parasite populations should accumulate these new gene variants when they are produced at a sufficient rate for their lifespans to overlap with each other (S1 Fig). Novelty per se guarantees neither the establishment nor the persistence of the genes. Even under high absolute fitness, new variants need to survive initial drift to establish in the parasite population [20]. Their accumulation further requires that the rate at which they are generated, G_{new} , be on average larger than that of their loss, given by the inverse of their lifespan, T_{new} . In other words, at least one beneficial gene needs to be produced and become established in the population during the typical lifespan of a previously generated new gene. We denote by R_{div} the expected number of new genes produced during the average lifespan. This reproductive or innovation number should be greater than 1 for new variants to accumulate, namely

$$R_{div} = G_{new}T_{new} = N\mu p_{inv}T_{new} > 1 \quad (1)$$

where N denotes the population size of the parasite, μ , the mutation rate of the genes, and p_{inv} , the invasion probability of a low-frequency variant. Importantly, the average lifespan of a new gene T_{new} is under frequency-dependent selection from hosts’

acquisition of immunity, with immune selection conferring an advantage to the rare and a disadvantage to the common [11,21]. New variants are generated through either mutation or ectopic recombination (Methods), and therefore μ generically refers to the rate of novelty generation regardless of specific mechanisms.

Equation (1) establishes an expectation for the existence of a threshold, whose expression we proceeded to further develop next, to be able to verify it computationally. The factors that determine the spread and establishment of a new antigenic variant consist of its selective advantage and the rates of innovation and transmission. We first consider the probability p_{inv} that a new gene survives its initial low frequency and invades. Based on birth-death processes in the Moran model with selection [22], the establishment probability of a low-frequency variant is largely determined by its relative fitness advantage over other genes. The fitness of genes in a transmission model is essentially given by their effective reproductive number R_{eff} or the number of copies they produce via transmission events during their lifetime. Thus,

$$p_{inv} \approx (R_{new}/\bar{R}_{eff}) - 1 \quad (2)$$

Equation (2) holds in general for any infectious disease that generates new antigens. To proceed further, we considered the specifics of *Plasmodium falciparum* and its multicopy *var* genes (typically about 40-60 per genome), whose expression is sequential during the blood stage of infection [23]. R_{eff} for a given *var* gene is the product of the epidemiological contact rate of the disease (β) and the typical infection duration (τ) of parasites that carry the given gene (Methods). If we consider that genes are equivalent in transmissibility (i.e., their products are functionally equivalent in their ability to bind host receptors), an assumption we later relax in numerical results, fitness differences between variants are only determined by the duration of infection these genes can typically sustain.

Because only those genes towards which the host has not yet built immunity are expressed, the average duration of infection will equal the number of genes per genome times the average proportion of susceptible (non-immune) hosts per gene,

$$\bar{\tau} = dg \sum_{i=1}^k S_i f_i \quad (3)$$

where d is the duration of infection for a given gene in a naive host, g is the number of genes per genome, and f_i is the population frequency of a given gene. We rewrote equation (2) using (3) (Methods), to obtain

$$p_{inv} \approx \frac{S_{new} - \bar{S}}{\bar{S}g} \quad (4)$$

where the mean number of susceptible hosts for a gene is given by $\bar{S} = \sum_{i=1}^k S_i f_i$. This expression for the invasion probability shows that a new gene is likely to invade when it affords a wider host niche than that of older genes (by encoding for epitopes that are new given the immunity of the host population). In other words, the available number of hosts for its expression should be higher on average than that for existing genes. In addition, the invasion probability of a single gene decreases with increasing genome size g , as the importance of a single gene also decreases.

We can now evaluate the existence of the threshold behavior indicated by the above analytical argument. To this end, we computed R_{div} from a stochastic agent-based model of malaria transmission [11,24] that tracks *var* evolution and immunity and is described in detail in [11]. A number of extensions were also considered here to address the generality of the argument, including distinct major *var* gene groups with associated differential fitness and constrained recombination (Methods). We specifically

examined how the accumulation of new variants over a given period of time varies as a function of R_{div} . We calculated G_{new} according to Eq. (2) and (4) by obtaining \bar{S} , N , and μ directly from the simulations, and T_{new} from the average lifespan of all the new variants that are produced during this time period.

Results showed that the transmission system naturally falls into two regimes separated by a threshold at $R_{div} = 1$ (Fig. 1A). Below this transition, new antigenic variants are generated but do not accumulate or persist (Fig. 1C), whereas above it, they are able to accumulate and experience a continuous turnover rate (see shifting shades of colors in Fig. 1D). The transition between these regimes occurs around the proposed boundary where the rate of generation of genes surpasses the average lifespan of new antigenic variants ($R_{div} > 1$). This threshold is robust to differences in specific assumptions about the transmission and genetic systems (including processes of within-host dynamics, functional differences between genes, values of the recombination and biting rates), as each point in Fig. 1A represents a simulation with different model assumptions and parameter combinations (Methods; S1-2 Table).

Importantly, we found that the quantity R_{div} scales monotonically with the intensity of transmission measured here as the entomological inoculation rate (or EIR , the number of infectious bites per person per year) (S2 Fig, Fig. 1B), a practical empirical measure from field epidemiology. The association with R_{div} should hold more generally with any other measure of transmission intensity. This association implies that the transition between regimes also occurs as a function of transmission intensity (Fig. 1A, C, D), and therefore, that the malaria system can be pushed below threshold by changing this control variable.

The transition examined so far represents the behavior of the system for different values of R_{div} or transmission intensity. Its existence should influence the temporal response of the malaria system to intervention events that reduce transmission at a given point in time. In particular, interventions that take the transmission system above threshold should lead to distinct responses than those that fail to do so. This is illustrated in Fig. 1 (G and H), where we numerically introduced a transient reduction of the biting rate to lower levels (respectively 30 and 50% of its original value). New genes cease to accumulate only when R_{div} goes below threshold (Fig. 1G), whereas they continue to invade and accumulate following a temporary decrease otherwise (Fig. 1H).

In order to further understand how epidemiological and genetic factors influence R_{div} , we examined with a simple theoretical model the equilibrium values of \bar{S} and N , which enter prominently in the expression for G_{new} (Methods). At equilibrium, increasing contact rates (β) via mosquito bites, result in higher parasite population sizes (N) and a lower average number of susceptible hosts (\bar{S} , S2 Fig). More specifically, \bar{S} is mostly determined by β and g , whereas parasite population size strongly scales with the diversity ratio (i.e., gene pool size G over genome size, G/g). A lower \bar{S} favors invasion of a new variant (i.e., increases p_{inv}), and a higher parasite population size (N) and a genome with a higher number of unique genes (g) generate new variants faster. Theoretical predictions underestimate \bar{S} because they neglect the higher-level of organization of the genes into different genomes, also under immune selection [11]. This indicates that such strain structure can significantly reduce the percentage of genes that a host is immune to, especially under high competition and high diversity (see $g = 60$, and $G/g = 100$ in S3B Fig).

An explicit expression for T_{new} cannot be obtained analytically because this quantity continuously changes as new genes enter the system and influence the nonlinear dynamics of N , G , and \bar{S} . To gain nevertheless an understanding for how these variables affect T_{new} , we approximated the average lifespan of a new gene on the basis of an adapted diffusion equation [22] under the assumption that the system remains constant and only this lifespan, we call \bar{t} , varies (SI). The diffusion equation for

\bar{t} requires consideration of how the frequency of a new gene $x(t)$ varies in time (SI).
When applied to our model, the resulting analytical approximation for \bar{t} (Eq. S8-11)
showed that the expected lifespan of a new gene grows faster than exponential with
decreasing \bar{S} , and surpasses the average time to fixation of a neutral gene ($2N$) when \bar{S}
is below a given value (40%) (Fig. 2). The mean T_{new} evaluated numerically will
always be shorter than that predicted from the diffusion approximation \bar{t} , especially as
 \bar{S} becomes smaller and persistence times rapidly increase. This is because the stochastic
simulations can only track lifespan within a finite time period (which places an upper
bound on its value), and because the assumption of constant \bar{S} does not apply.
Nevertheless, the theoretical trend of increasing T_{new} with decreasing \bar{S} and therefore,
higher transmission intensity, does apply to the numerical system.

In summary, by evaluating whether $G_{new} < 1/T_{new}$ (or equivalently, $R_{div} < 1$), one
can predict whether the system has a relatively stable antigen-encoding gene pool, or
whether alternatively, new variants continuously enter into it. We have shown that
whether new genes can successfully establish in the population is most tightly linked
with the average proportion of susceptible hosts (or the niche) available for existing
genes (\bar{S}). When transmission rate is low, \bar{S} is large and new genes do not have a
significant advantage over older ones. New genes experience a small invasion probability,
and even when they invade, they experience strong drift, functioning as effectively
neutral. As transmission intensity increases, the selective advantage of new genes also
increases as \bar{S} decreases. Once \bar{S} is below a given value (0.4 in our simulations), new
genes are most likely to be maintained in the population indefinitely. Concomitantly,
the increase in gene diversity results in higher parasite population sizes N (S3A Fig).
The system thus enters a regime of positive feedback for new variants, as elevated
diversity boosts N and therefore also, G_{new} , before reaching equilibrium.

Discussion

The concept of R_{div} arises from the interplay of immune memory and antigenic
variation at the population level, as a result of frequency-dependent selection. As such,
it differs from the antigenic diversity threshold previously proposed for the HIV virus
and its transition to AIDS, arising from the race between viral replication and immune
responses at the within-host level [25]. The concept itself and the associated transition
regime described here should apply more generally to other infectious diseases with
antigen-encoding multigene families, such as *vsg* genes in *Trypanosoma brucei* and *msg*
genes in *Pneumocystis carinii* [18]. Because the basic concept is independent of specific
consideration of multigene families and their properties, it should also be adaptable to
other pathogens in which large standing antigenic diversity at the population level
results from multilocus genetic variation [5]. By contrast in pathogens with sufficiently
well-defined strains, R_0 would be sufficient because the characteristics of their
population dynamics and population genetics would keep them below the diversification
threshold defined by R_{div} equal one. For example, genetic variation in measles is largely
neutral antigenically and the effective mutation rates generating new antigens are
slow [26]. In seasonal influenza, bottlenecks in transmission constrain the emergence of
novelty [27], and so do mutations with largely deleterious effects [28].

For falciparum malaria and pathogens with extensive antigenic diversity at the
population level, the proposed concept of a threshold behavior in the accumulation of
antigen-encoding genes has practical implications for overcoming the resilience of highly
endemic regions to intervention efforts. Although a decreasing trend in the diversity of
strains and underlying genes with transmission intensity is well known and expected
from both the biogeography and epidemiology of malaria, the actual form of this
reduction is much less clear. Our results predict the existence of a sharp transition

below which the disease system should effectively respond as a typical low transmission region, not just because of reduced transmission intensity but also because of much lower antigenic diversity no longer able to rebuild. Failure to push transmission intensity below this threshold would lead to a fast rebound in new antigenic variation, despite an overall diversity reduction. The crossing of the threshold would instead provide an indication that the system is now poised for further intervention with enhanced results.

Control and even elimination efforts are indeed known to be most successful in biogeographical regions of low transmission, such as those at the edge of the distribution of the disease in Africa and in other continents [29]. Arresting the fast turnover of the local antigenic pool typical of high endemism would significantly repress disease burden and facilitate its further reduction. Concomitant control efforts at a regional level are critical to stem immigration, as migrant genomes would exhibit higher invasion probabilities than local ones, given their higher likelihood of harboring new antigens. Monitoring the turnover of *var* gene diversity through molecular epidemiology in response to control efforts should inform intervention evaluation in high transmission regions.

Although the importance of host immune selection in shaping the antigenic variation of *P. falciparum* and other pathogens is recognized [21, 30–32], mathematical and computational models typically evaluate intervention efficacy without explicit consideration of antigenic diversity (e.g. [33, 34]) and openness of the system to innovation [35]. Our results underscore the importance of these aspects.

Estimation of the diversification number, R_{div} , would provide general guidelines for intervention evaluation where traditional application of R_0 is unable to do so for highly diverse pathogens [2, 21]. Future work should consider how to obtain this number from an estimation of key parameters, including parasite population size, transmission rates, and gene pool size, based on combined data from molecular and field epidemiology. Parameterization of an agent-based stochastic transmission model that implements immune selection and recombination explicitly (e.g., [11, 36]) could be used, which represents a computational challenge (but see [37]). Estimating the viability of new recombinants will require bioinformatic analyses of population-level *var* sequence data for the DBL_α portion of the gene [38].

For simplicity, our analytical derivations treated each gene independently, even though *var* gene composition in parasite genomes of local populations in regions of high transmission has been shown to be non-randomly and non-neutrally structured, exhibiting low overlap as the result of immune selection [11, 13, 24, 39]. Hence, the fate of a viable new antigen-encoding gene depends on its genomic background, which ultimately determines the strength of competition among parasites. Comparisons of analytical expectations with numerical simulations revealed an influence of such population structure on the fate of new genes, and therefore, on components of R_{div} . Future work should examine extensions of this work that account for this further complexity of immune selection operating at different levels of organization.

In general, explicit consideration of a reproductive number for antigenic diversification should enhance our understanding of transmission dynamics where large standing pathogen diversity represents a major challenge to control efforts.

Materials and methods

Analytical derivation of R_{div}

We consider a population of hosts whose number is denoted by N_{host} , receiving malaria infections from a diverse set of parasites, each composed of g genes from a constant gene pool of size G . One of the two main components of R_{div} is the rate at which new genes

are generated, G_{new} . Besides the mutation rate μ and the equilibrium parasite population size N , its expression requires the invasion probability p_{inv} , we derive below.

Invasion probability of a new variant, p_{inv}

From birth-death processes according to the Moran model with selection [22], the probability of establishment of a low frequency variant is determined by its fitness advantage relative to that of other genes, and by the parasite population size. That is,

$$p_{inv} = \frac{1 - (W/W_{new})^n}{1 - (W/W_{new})^N} \quad (5)$$

where n denotes the number of copies of new genes, and W , the fitness of a gene. In our case, $n = 1$ as new genes originate from a unique mutation or an ectopic recombination event. When $N \gg 1$, the invasion probability p_{inv} is approximately its initial selective advantage relative to established gene variants, provided the selection coefficient remains the same. Since the fitness of each individual gene in a transmission model is essentially given by their effective reproductive number R_{eff} , we have

$$p_{inv} \approx (R_{new}/\bar{R}_{eff}) - 1. \quad (6)$$

R_{eff} for a given *var* gene is in turn the product of the epidemiological contact rate of the disease (β) and the typical infection duration (τ) of parasites that carry the given gene,

$$R_{eff} = \bar{\beta}\bar{\tau} \quad (7)$$

The contact rate β is equal to the product of the transmission rate (b) and the ‘transmissibility’ or infectivity of the given *var* gene (i.e., the functionality of the gene, f). Because we do not model vectors explicitly in the numerical stochastic model, the contact rate (β) refers to the rate at which a transmission event occurs, with a ‘donor’ host transmitting infection to a ‘recipient’ one (detailed description in section on “the modified *var* evolution model”).

Different groups of *var* genes may vary in their binding affinities to host receptors and therefore in their transmissibility. For simplicity, we consider that all genes exhibit the same transmissibility and therefore, the same absolute fitness, as we are most interested in estimating the fate of a new variant as a result of immune (frequency-dependent) selection. (We do explore later the effect of fitness differences numerically with the agent-based stochastic model, described in the section on “the modified *var* evolution model”).

With (7) for R_{eff} , we can write

$$\begin{aligned} p_{inv} &\approx \frac{\tau_{new}}{\bar{\tau}} - 1 \\ &= \frac{\bar{S} \times (g - 1) + S_{new}}{\bar{S} \times g} - 1 \end{aligned} \quad (8)$$

Numerical evaluation of R_{div}

With equation (8) (or its equivalent (4)), we can now compute $G_{new} = N\mu p_{inv}$ in equation (1) from the output of our numerical simulations (described below under the *The modified var evolution model*). We also obtained from the simulations the other major component of R_{div} , the mean lifetime of the genes in the system, T_{new} , by directly tracking their fate individually. Although an analytical expression for T_{new} was not achievable for this nonlinear and stochastic transmission system, we considered gene lifetime under simplifying assumptions, as explained next.

Analytical derivation of \bar{t}

The simplifying assumptions are that the system has reached an equilibrium ((i.e., parasites get transmitted and die at the same rate), and that only the average lifetime of a gene varies, with all other variables remaining unchanged, including N and the average proportion of hosts \bar{S} susceptible to an average gene. To differentiate gene lifetime under these conditions from T_{new} itself, we call it \bar{t} . An expression for \bar{t} is derived by considering the frequency-dependent selection experienced by a new gene variant entering the system at equilibrium. We specifically approximate the dynamics of \bar{t} on the basis of an adapted diffusion equation [22] (Supplementary Text).

The modified *var* evolution model

We used an extended implementation of the agent-based model developed in [11], where a complete description can be found. Here, we first briefly summarize the main features of the computational model, and then document the specific changes implemented in this study, including different transmission scenarios and rules of within-host dynamics. (Parameter combinations and specific rules are listed in Table S1-2).

The computational model is an individual-based, discrete-event, continuous-time stochastic system in which the infection and immune history of each host are tracked individually. In the numerical implementation of the simulation, all possible future events are stored in a single event queue along with their putative times, which may be fixed or drawn from a probability distribution. When an event occurs, it may trigger the addition or removal of future events on the queue, or changes of their rates, leading to a recalculation of their putative time. The implementation is adapted from the next-reaction method [40], which optimizes the Gillespie first-reaction method [41] and allows for faster simulation times.

Transmission events are sampled at the rate, $N_{host}\beta$, in which a donor and a recipient host are sampled randomly from the host population. If the donor harbors parasites, then each parasite has a probability of being transmitted to the mosquito that is proportional to the functionality of the *var* gene that is currently under expression. *var* genomes picked up by the mosquito will recombine with another genome to produce sporozoites. Specifically, if there are n parasite genomes, each genome has a probability $1/n$ of recombining with itself, producing the same offspring genome, and a probability $1 - 1/n$ of recombining with a different genome, producing recombinants. The total number of *var* genomes passed onto the receiver is kept the same as that received from biting the infectious host. Parasites in a human host are not infectious until the completion of a delay, representing altogether oocysts development in the mosquito during the sexual stage and the initial liver stage in the receiver host. Since we do not model mosquitoes explicitly, we implement this delay as a 14-day gap between a transmission event and the genome becoming infectious. *var* genes within a genome express sequentially in a random order, or according to their functional level from high to low, depending on the specific rule (Table S1), with a switching rate of $1/d$ (with d denoting the mean duration of expression). When the expression switches to a new gene, immunity to the gene previously expressed is added to the host's immune memory. The infection ends when the expression of the whole *var* genome is completed. Mutation μ and ectopic recombination r occur randomly during the infectious stage (see below).

Genome structure

The genome of an individual parasite is a combination of g *var* genes. Each *var* gene, in turn, is a linear combination of two loci encoding epitopes that are connected linearly, and each epitope can be viewed as a multi-allele locus with n possible alleles. The

initial conditions for the simulation include g^*2 alleles per epitope and g^*20 combinations of these genes in the gene pool. A typical simulation starts by initializing the local parasite population via a given number of transmission events with migrant genomes whose composition is sampled randomly from a regional pool of genes, G . Specifically, 20 hosts are infected with randomly assembled genomes, and one migrant genome is introduced every day into the population to simulate exposure to all the variants in the genome pool quickly.

Ectopic recombination during the asexual blood stage of infection

var genes often change their physical location and form new variants through ectopic recombination and gene conversions. These processes occur during both sexual and asexual stages. As ectopic recombination is observed more often during the asexual stage where parasites spend most of their life cycle, and our model does not represent the mosquito stage explicitly, we consider ectopic recombination among genes within the same genome during the asexual stage. Two genes are selected from the genome repertoire, with a breakpoint located along the gene randomly. Newly recombined genes have a probability P_f to be functional (i.e., viable), defined by the similarity of the new variant with its parental genes as in He et al. [11]. If the recombined gene is selected to be non-functional, then the parental gene is kept. Otherwise, the recombined gene substitutes the parental one and a new strain is formed. In the current implementation, each recombination has a 50% chance to generate a new allele.

var gene groups and trade-offs

In the previous model, genes differed antigenically but not functionally. For increased realism, each gene is assigned to either *var* upsGroup A or *var* upsGroup B/C to represent existing differences in recombination rates and functionality of *var* gene groups [42]. Ectopic recombination is only allowed to occur within each group, and genes in upsGroup B/C have higher recombination rates than those in upsGroup A [6]. Each gene is also assigned an intrinsic growth rate of the parasites f that express it (Table S2), because antigens that better bind to host receptors result in a higher parasite growth rate [43,44]. In an additional implementation of the model, genes with higher growth rates are expressed first, followed in decreasing order by genes of lower growth rates. Also, genes with higher growth rates are expected to be cleared faster by the immune system, translating into a higher switch rate to the next gene, which is controlled by the trade-off parameter tff .

Data Availability

The original C++ code for the *var* evolution model is available on Github (<https://github.com/pascualgroup/VarModel2>).

Supporting information

S1 Appendix. Supplementary text for theoretical derivations.

S1 Fig. Schematic illustration of two possible scenarios.

S2 Fig. Relationship between transmission intensity and R_{div} .

S3 Fig. Comparison between theoretical expectations from Eq. S1 (\diamond) and corresponding values from stochastic simulations (\bullet) for N (A) and \bar{S} (B) as a function of contact rate, β , genome size, g , and two levels of diversity ratio, $G/g = 10$ or 100 . 383
384
385
386

S4 Fig. The deterministic trajectory of a new gene variant invading a system that is previously at equilibrium under a low (0.015, left panels) and a high (0.05, right panels) contact rate. 387
388
389

S5 Fig. Phase diagram of $x(t)$ and $S_{new}(t)$ from S4 Fig. 390

S6 Fig. Persistence of new genes according to \bar{t} . 391

S1 Table. Epidemiological and genetic parameters used in stochastic simulations. 392
393

S2 Table. Epidemiological, genetic and within-host dynamics rules varied in the stochastic simulations. 394
395

Acknowledgments 396

We are grateful to the funding provided by the joint NIH-NSF-NIFA Ecology and Evolution of Infectious Disease (award R01 AI149779). We acknowledge valuable discussions at the Santa Fe Institute, as part of a working group supported by the James S. McDonnell Foundation (JSMF). We thank Karen P. Day for her insightful comments on an earlier version of the manuscript. We appreciate the support of the University of Chicago through computational resources at the Midway cluster. 397
398
399
400
401
402

References

1. Keeling MJ, Rohani P. Modeling Infectious Diseases in Humans and Animals. Princeton University Press; 2008.
2. Gupta S, P Day K. A theoretical framework for the immunoepidemiology of Plasmodium falciparum malaria. Parasite Immunology. 1994;16(7):361–370. doi:10.1111/j.1365-3024.1994.tb00361.x.
3. Buckee CO, Recker M. Evolution of the Multi-Domain Structures of Virulence Genes in the Human Malaria Parasite, Plasmodium falciparum. PLOS Computational Biology. 2012;8(4):e1002451. doi:10.1371/journal.pcbi.1002451.
4. Deitsch KW, Moxon ER, Wellems TE. Shared themes of antigenic variation and virulence in bacterial, protozoal, and fungal infections. Microbiol Mol Biol Rev. 1997;61(3):281–293.
5. Georgieva M, Buckee CO, Lipsitch M. Models of immune selection for multi-locus antigenic diversity of pathogens. Nature Reviews Immunology. 2019;19(1):55. doi:10.1038/s41577-018-0092-5.
6. Claessens A, Hamilton WL, Kekre M, Otto TD, Faizullahoy A, Rayner JC, et al. Generation of Antigenic Diversity in Plasmodium falciparum by Structured Rearrangement of Var Genes During Mitosis. PLOS Genetics. 2014;10(12):e1004812. doi:10.1371/journal.pgen.1004812.

7. Frank M, Kirkman L, Costantini D, Sanyal S, Lavazec C, Templeton TJ, et al. Frequent recombination events generate diversity within the multi-copy variant antigen gene families of *Plasmodium falciparum*. *International Journal for Parasitology*. 2008;38(10):1099–1109. doi:10.1016/j.ijpara.2008.01.010.
8. Chen DS, Barry AE, Leliwa-Sytek A, Smith TA, Peterson I, Brown SM, et al. A Molecular Epidemiological Study of var Gene Diversity to Characterize the Reservoir of *Plasmodium falciparum* in Humans in Africa. *PLOS ONE*. 2011;6(2):e16629. doi:10.1371/journal.pone.0016629.
9. Day KP, Artzy-Randrup Y, Tiedje KE, Rougeron V, Chen DS, Rask TS, et al. Evidence of strain structure in *Plasmodium falciparum* var gene repertoires in children from Gabon, West Africa. *Proceedings of the National Academy of Sciences*. 2017;114(20):E4103–E4111. doi:10.1073/pnas.1613018114.
10. Ruybal-Pesántez S, Tiedje KE, Tonkin-Hill G, Rask TS, Kanya MR, Greenhouse B, et al. Population genomics of virulence genes of *Plasmodium falciparum* in clinical isolates from Uganda. *Scientific Reports*. 2017;7(1):11810. doi:10.1038/s41598-017-11814-9.
11. He Q, Pilosof S, Tiedje KE, Ruybal-Pesántez S, Artzy-Randrup Y, Baskerville EB, et al. Networks of genetic similarity reveal non-neutral processes shape strain structure in *Plasmodium falciparum*. *Nature Communications*. 2018;9(1):1817. doi:10.1038/s41467-018-04219-3.
12. Ruybal-Pesántez S, Tiedje KE, Pilosof S, Tonkin-Hill G, He Q, Rask TS, et al. Why do adults harbor *Plasmodium falciparum* infections in high transmission African settings? *eLife*. 2020;.
13. Day KP, Artzy-Randrup Y, Tiedje KE, Rougeron V, Chen DS, Rask TS, et al. Evidence of strain structure in *Plasmodium falciparum* var gene repertoires in children from Gabon, West Africa. *PNAS*. 2017;114(20):E4103–E4111. doi:10.1073/pnas.1613018114.
14. Ruybal-Pesántez S, Tiedje KE, Tonkin-Hill G, Rask TS, Kanya MR, Greenhouse B, et al. Population genomics of virulence genes of *Plasmodium falciparum* in clinical isolates from Uganda. *Sci Rep*. 2017;7(1):11810. doi:10.1038/s41598-017-11814-9.
15. Barnwell JW, Howard RJ, Coon HG, Miller LH. Splenic requirement for antigenic variation and expression of the variant antigen on the erythrocyte membrane in cloned *Plasmodium knowlesi* malaria. *Infection and Immunity*. 1983;40(3):985–994.
16. Handunnetti SM, Mendis KN, David PH. Antigenic variation of cloned *Plasmodium fragile* in its natural host *Macaca sinica*. Sequential appearance of successive variant antigenic types. *The Journal of Experimental Medicine*. 1987;165(5):1269–1283. doi:10.1084/jem.165.5.1269.
17. Stern A, Brown M, Nickel P, Meyer TF. Opacity genes in *Neisseria gonorrhoeae*: Control of phase and antigenic variation. *Cell*. 1986;47(1):61–71. doi:10.1016/0092-8674(86)90366-1.
18. Deitsch KW, Lukehart SA, Stringer JR. Common strategies for antigenic variation by bacterial, fungal and protozoan pathogens. *Nat Rev Micro*. 2009;7(7):493–503. doi:10.1038/nrmicro2145.

19. Levine JM, Bascompte J, Adler PB, Allesina S. Beyond pairwise mechanisms of species coexistence in complex communities. *Nature*. 2017;546(7656):56–64. doi:10.1038/nature22898.
20. Hermisson J, Pennings PS. Soft Sweeps: Molecular Population Genetics of Adaptation From Standing Genetic Variation. *Genetics*. 2005;169(4):2335–2352. doi:10.1534/genetics.104.036947.
21. Gupta S, Maiden MCJ, Feavers IM, Nee S, May RM, Anderson RM. The maintenance of strain structure in populations of recombining infectious agents. *Nat Med*. 1996;2(4):437–442. doi:10.1038/nm0496-437.
22. Ewens WJ. *Mathematical Population Genetics, I. Theoretical Introduction*. 2; 2004.
23. Recker M, Buckee CO, Serazin A, Kyes S, Pinches R, Christodoulou Z, et al. Antigenic Variation in *Plasmodium falciparum* Malaria Involves a Highly Structured Switching Pattern. *PLOS Pathogens*. 2011;7(3):e1001306. doi:10.1371/journal.ppat.1001306.
24. Artzy-Randrup Y, Rorick MM, Day K, Chen D, Dobson AP, Pascual M. Population structuring of multi-copy, antigen-encoding genes in *Plasmodium falciparum*. *Elife*. 2012;1:e00093.
25. Nowak MA, May RM. Mathematical biology of HIV infections: antigenic variation and diversity threshold. *Mathematical Biosciences*. 1991;106(1):1–21. doi:10.1016/0025-5564(91)90037-J.
26. Beaty SM, Lee B. Constraints on the Genetic and Antigenic Variability of Measles Virus. *Viruses*. 2016;8(4):109. doi:10.3390/v8040109.
27. McCrone JT, Woods RJ, Martin ET, Malosh RE, Monto AS, Lauring AS. Stochastic processes constrain the within and between host evolution of influenza virus. *eLife*. 2018;7:e35962. doi:10.7554/eLife.35962.
28. Fitzsimmons WJ, Woods RJ, McCrone JT, Woodman A, Arnold JJ, Yennawar M, et al. A speed–fidelity trade-off determines the mutation rate and virulence of an RNA virus. *PLOS Biology*. 2018;16(6):e2006459. doi:10.1371/journal.pbio.2006459.
29. Hemingway J, Shretta R, Wells TNC, Bell D, Djimdé AA, Achee N, et al. Tools and Strategies for Malaria Control and Elimination: What Do We Need to Achieve a Grand Convergence in Malaria? *PLOS Biology*. 2016;14(3):e1002380. doi:10.1371/journal.pbio.1002380.
30. Buckee CO, Jolley KA, Recker M, Penman B, Kriz P, Gupta S, et al. Role of selection in the emergence of lineages and the evolution of virulence in *Neisseria meningitidis*. *PNAS*. 2008;105(39):15082–15087. doi:10.1073/pnas.0712019105.
31. Gupta S, Ferguson N, Anderson R. Chaos, Persistence, and Evolution of Strain Structure in Antigenically Diverse Infectious Agents. *Science*. 1998;280(5365):912–915. doi:10.1126/science.280.5365.912.
32. Recker M, Nee S, Bull PC, Kinyanjui S, Marsh K, Newbold C, et al. Transient cross-reactive immune responses can orchestrate antigenic variation in malaria. *Nature*. 2004;429(6991):555. doi:10.1038/nature02486.

33. Griffin JT, Hollingsworth TD, Okell LC, Churcher TS, White M, Hinsley W, et al. Reducing *Plasmodium falciparum* Malaria Transmission in Africa: A Model-Based Evaluation of Intervention Strategies. *PLOS Medicine*. 2010;7(8):e1000324. doi:10.1371/journal.pmed.1000324.
34. Nkuo Akenji TK, Ntonifor NN, Ching JK, Kimbi HK, Ndamukong KN, Anong DN, et al. Evaluating a malaria intervention strategy using knowledge, practices and coverage surveys in rural Bolifamba, southwest Cameroon. *Transactions of The Royal Society of Tropical Medicine and Hygiene*. 2005;99(5):325–332. doi:10.1016/j.trstmh.2003.12.016.
35. Holding T, Valletta JJ, Recker M. Multiscale Immune Selection and the Transmission-Diversity Feedback in Antigenically Diverse Pathogen Systems. *The American Naturalist*. 2018;192(6):E189–E201. doi:10.1086/699535.
36. Buckee CO, Recker M, Watkins ER, Gupta S. Role of stochastic processes in maintaining discrete strain structure in antigenically diverse pathogen populations. *PNAS*. 2011;108(37):15504–15509. doi:10.1073/pnas.1102445108.
37. Ozik J, Collier NT, Wozniak JM, Spagnuolo C. From desktop to Large-Scale Model Exploration with Swift/T. In: 2016 Winter Simulation Conference (WSC); 2016. p. 206–220.
38. Zilversmit MM, Chase EK, Chen DS, Awadalla P, Day KP, McVean G. Hypervariable antigen genes in malaria have ancient roots. *BMC Evolutionary Biology*. 2013;13:110. doi:10.1186/1471-2148-13-110.
39. Pilosof S, He Q, Tiedje KE, Ruybal-Pesántez S, Day KP, Pascual M. Competition for hosts modulates vast antigenic diversity to generate persistent strain structure in *Plasmodium falciparum*. *PLOS Biology*. 2019;17(6):e3000336. doi:10.1371/journal.pbio.3000336.
40. Gibson MA, Bruck J. Efficient exact stochastic simulation of chemical systems with many species and many channels. *J Phys Chem A*. 2000;104:1876—1889.
41. Gillespie DT. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics*. 1976;22(4):403–434. doi:10.1016/0021-9991(76)90041-3.
42. Rowe JA, Claessens A, Corrigan RA, Arman M. Adhesion of *Plasmodium falciparum*-infected erythrocytes to human cells: molecular mechanisms and therapeutic implications. *Expert Reviews in Molecular Medicine*. 2009;11. doi:10.1017/S1462399409001082.
43. Bachmann A, Petter M, Krumkamp R, Esen M, Held J, Scholz JAM, et al. Mosquito Passage Dramatically Changes var Gene Expression in Controlled Human *Plasmodium falciparum* Infections. *PLOS Pathogens*. 2016;12(4):e1005538. doi:10.1371/journal.ppat.1005538.
44. Bachmann A, Bruske E, Krumkamp R, Turner L, Wichers JS, Petter M, et al. Controlled human malaria infection with *Plasmodium falciparum* demonstrates impact of naturally acquired immunity on virulence gene expression. *PLOS Pathogens*. 2019;15(7):e1007906. doi:10.1371/journal.ppat.1007906.

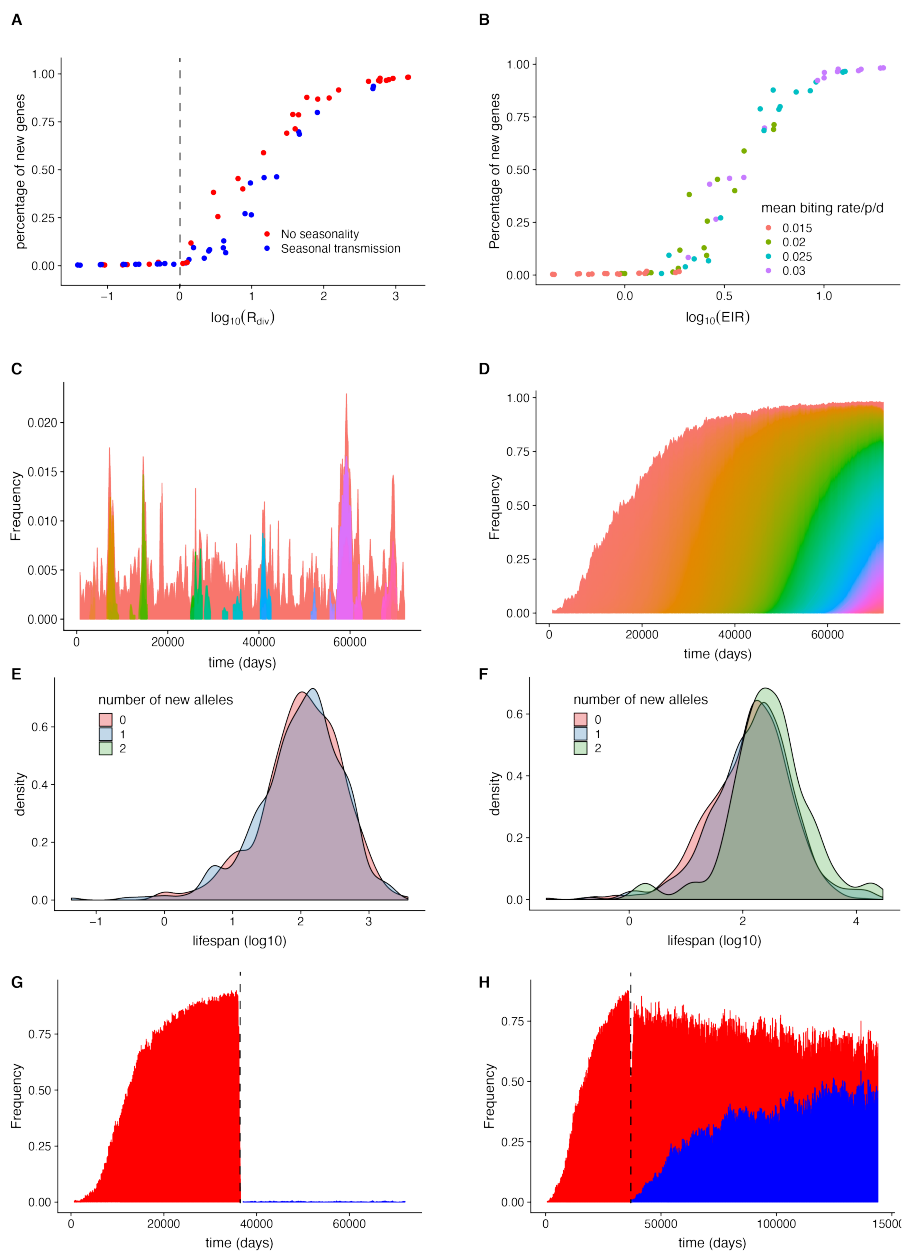


Fig 1. Numerical simulations reveal a transition between two regimes of antigenic diversity accumulation. (A) The percentage of new genes in the local parasite population at the end of a given simulation period (200 years) remains negligible when the reproductive number R_{div} for antigen-encoding new genes is lower than one. By contrast, this percentage increases rapidly above this threshold. Because the time interval over which we computed $R_{div} = G_{new}T_{new}$ concerns long transients, we evaluated the rate of generation of new genes G_{new} as a mean over this interval (by averaging the values of N and \bar{S} every 180-day interval), and the lifespan T_{new} , as an average for all the new genes that invaded during this time (with this interval placing an upper bound on individual lifespans). Each point represents a simulation with different combinations of parameters and assumptions (including variation in rules of within-host dynamics, in strength of the trade-off between transmissibility and duration of infection, and in values and seasonality of the transmission rates, Table S1-2). (B) Because R_{div} increases monotonically with transmission intensity (Fig. S2), the percentage of new genes also exhibits the threshold behavior with this variable, measured here by the entomological inoculation rate (EIR, the number of infectious bites per person per year). For simplicity, when a transmission event occurs, our model considers that all bites of an infected ‘donor’ host generate an infectious bite, and that all infectious bites of the ‘recipient’ host result in infection. This implies that the EIR values in the graph should be adjusted for comparison to actual field values (by dividing by the product of the competence/transmissibility probabilities, which will raise EIR). (C) New genes do not accumulate below the transmission threshold where they essentially follow neutral dynamics. In contrast, they do accumulate and turn over at a constant rate under frequency-dependent selection above this threshold (D). Each color in these two panels refers to a new gene in the population. New genes that account for less than five infections over the entire period are combined and represented in red. (E) The average lifespan of new genes is shorter below the threshold than above it (F). (F) In addition, new genes with a greater number of new alleles (epitopes) live longer above the threshold, whereas below the threshold, they experience similar lifespans. Interventions that push the system below the threshold are effective at stopping the accumulation of new genes (G), whereas those that do not, result in the rebound and rebuilding of diversity (H). Red and blue colors indicate genes that originate respectively before and after the intervention.

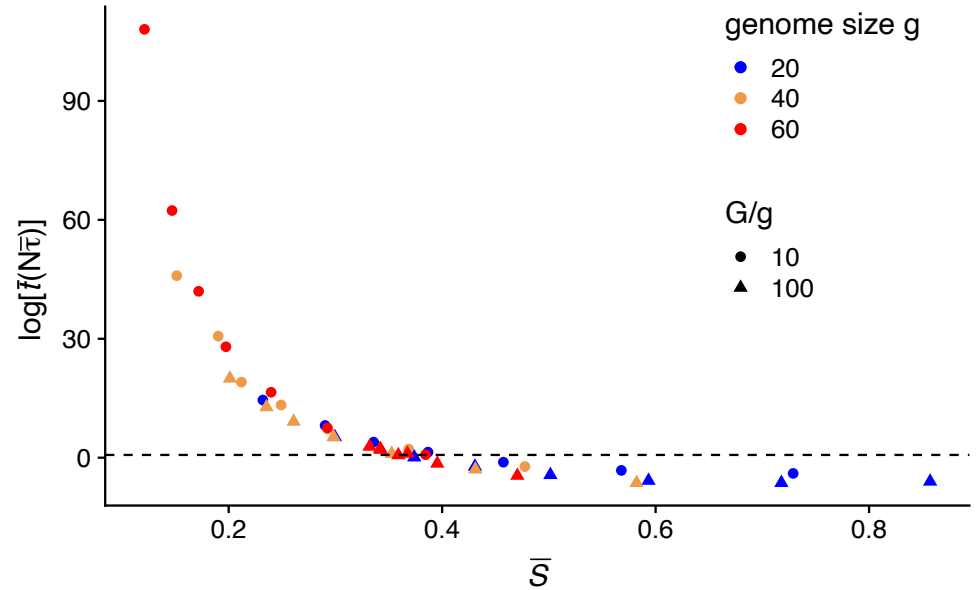


Fig 2. Theoretical expectation of the average lifespan of a new gene \bar{t} . The analytical expression shows that \bar{t} , measured in units of $N\bar{\tau}$, increases faster than exponential as the average number of available hosts \bar{S} decreases. The dashed line represents the time to fixation of a neutral gene, which means that under small \bar{S} , once established, the gene can be maintained in the population for much longer than the typical epidemiological timescale (or for much longer than the simulation period of 200 years in our model). The average lifespan T_{new} obtained from the computational model will always be considerably smaller than the theoretical expectation \bar{t} derived under the assumption that other factors remain constant, in particular the average number of hosts \bar{S} that are susceptible to the invading gene. The general trend of increased persistence with lower \bar{S} will hold however for the full numerical system and for finite time windows.