

Pathogen selection drives nonoverlapping associations between HLA loci

Bridget S. Penman^a, Ben Ashby^a, Caroline O. Buckee^b, and Sunetra Gupta^{a,1}

^aDepartment of Zoology, University of Oxford, Oxford OX1 3PS, United Kingdom; and ^bCenter for Communicable Disease Dynamics, Department of Epidemiology, Harvard School of Public Health, Boston, MA 02115

Edited* by Robert M. May, University of Oxford, Oxford, United Kingdom, and approved October 14, 2013 (received for review March 21, 2013)

Pathogen-mediated selection is commonly invoked as an explanation for the exceptional polymorphism of the HLA gene cluster, but its role in generating and maintaining linkage disequilibrium between HLA loci is unclear. Here we show that pathogen-mediated selection can promote nonrandom associations between HLA loci. These associations may be distinguished from linkage disequilibrium generated by other population genetic processes by virtue of being nonoverlapping as well as nonrandom. Within our framework, immune selection forces the pathogen population to exist as a set of antigenically discrete strains; this then drives nonoverlapping associations between the HLA loci through which recognition of these antigens is mediated. We demonstrate that this signature of pathogen-driven selection can be observed in existing data, and propose that analyses of HLA population structure can be combined with laboratory studies to help us uncover the functional relationships between HLA alleles. In a wider coevolutionary context, our framework also shows that the inclusion of memory immunity can lead to robust cyclical dynamics across a range of host–pathogen systems.

infectious disease | major histocompatibility complex | mathematical model | human evolution | population genetics

HLA, found on the surface of all nucleated cells, present pathogen peptides to T lymphocytes and are thus a keystone of adaptive immunity. Demonstrable associations of particular HLA alleles with resistance or susceptibility to severe disease (1, 2) underscore the importance of their role in protection against death from infection. The genes encoding HLAs are found in the 3.6-Mb-long MHC on chromosome 6 and are distinguished by their exceptional polymorphism (3), which is likely the result of selection from pathogens (4–6). Despite this enormous diversity, most human populations are dominated by a relatively small number of combinations of the alleles present at the class I HLA (A, B, C) and the principal class II HLA (DP, DQ, and DR) loci (7–12). Here we present a coevolutionary model demonstrating that pathogen selection can drive such long-term, long-range associations between HLAs. We show that this mechanistic process can be distinguished from other evolutionary effects by virtue of generating a higher degree of nonoverlap between HLA repertoires than might be expected under founder effects or hitchhiking.

A Multilocus Model for Host–Pathogen Coevolution with Allele-Specific Adaptive Immunity

We first explored the properties of a deterministic epidemiological model (*Methods*) in which (i) the pathogen population was represented by four potential strains defined by two antigenic loci containing alleles (a, b) and (x, y), respectively, and (ii) we defined within a diploid host, alleles (A,B) and (X,Y) at two linked “recognition loci” (i.e., HLA loci), each only capable of responding to the corresponding parasite allele (or epitope) given in lowercase above. We assumed that immunity developed in an allele-specific manner conferring complete protection against infection by any other antigenic type containing that

allele, but that there was a risk of death if a host was incapable of recognizing either allele of the infecting pathogen strain (Fig. 1).

In line with previous observations, the pathogen population was observed to adopt a discrete, nonoverlapping strain structure (13). However, once any two strains (e.g., ax and by) achieve dominance, host homozygotes AX/AX and BY/BY suffer from increased mortality because each is only able to mount an immune response against one of the two circulating pathogen strains (all other host genotypes can recognize at least one allele of both strains). The numbers of these homozygotes fall until eventually the only host haplotypes left in the population are AY and BX. Thus, the strain structuring of the pathogen population by host immune selection generates nonrandom associations among the immune recognition genes of the host.

This scenario will be stable (Fig. 2A) in the absence of pathogen mutation, or when the basic reproduction number of the pathogen (R_0 ; a measure of its fundamental transmission potential) (14) is low. Conversely, if R_0 is above a certain threshold, no genetic structuring is possible in either pathogen or host (*SI Appendix, Fig. S1*). Between these two extremes, we observe coevolutionary cycling (Fig. 2B) in place of permanent structuring. This dynamic emerges due to the fact that as soon as the pathogen population becomes dominated by a particular set of strains (say ax and by), haplotypes that are incapable of recognizing any one of the dominant pathogen strains (i.e., BY and AX) start to go down in frequency, and haplotypes that can recognize both the dominant pathogen strains (i.e., BX and AY) increase in frequency. Eventually the proportion of BX and AY in the population will be so high that it will be in the pathogen’s interest to switch its strain structure to (ay, bx) so as to exploit the infection reservoirs created by homozygotes of these haplotypes (BX/BX cannot become immune to ay). The system is capable of generating nonrandom associations between recognition alleles

Significance

Human leukocyte antigens (HLA), first identified in tissue-matching for transplantation, play a critical role in immunity. HLAs are extraordinarily diverse, but certain sets of HLA genes are more likely to be found together than others. Here, we show that associations between HLA genes can arise through their coevolutionary interaction with pathogens. Technological advances are making it easier to determine HLA types, but DNA sequence alone cannot fully predict an HLA’s functional properties. Our work offers a new evolutionary approach to tackling this problem.

Author contributions: S.G. designed research; B.S.P. performed research; B.A. and C.O.B. contributed new reagents/analytic tools; B.S.P. analyzed data; and B.S.P. and S.G. wrote the paper.

The authors declare no conflict of interest.

*This Direct Submission article had a prearranged editor.

Freely available online through the PNAS open access option.

¹To whom correspondence should be addressed. E-mail: sunetra.gupta@zoo.ox.ac.uk.

This article contains supporting information online at www.pnas.org/lookup/suppl/doi:10.1073/pnas.1304218110/-DCSupplemental.

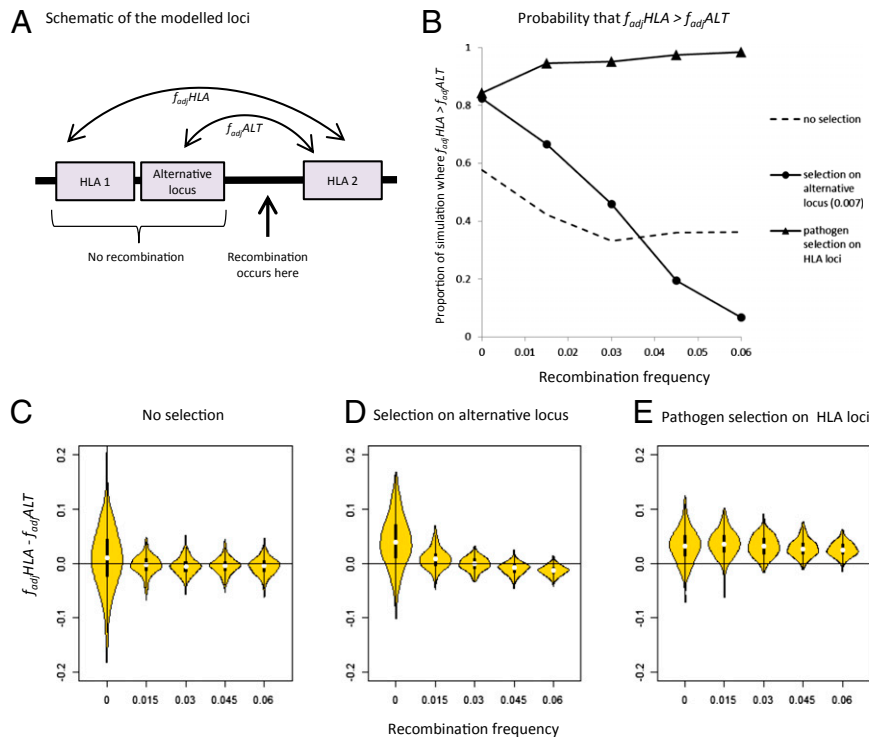


Fig. 5. Identifying a unique signature of HLA/pathogen coevolution. (A) Schematic representation of the loci within our model framework, indicating the pairs of loci between which $f_{\text{adj}}^{\text{HLA}}$ and $f_{\text{adj}}^{\text{ALT}}$ measure the degree of nonoverlap. (B) Probability that $f_{\text{adj}}^{\text{HLA}} > f_{\text{adj}}^{\text{ALT}}$ under different selection regimes and levels of recombination. (C–E) Violin plots (32) produced using the *vioplot* package in R version 2.15.2, representing the distribution of values of $f_{\text{adj}}^{\text{HLA}} - f_{\text{adj}}^{\text{ALT}}$ for 350 simulations at each indicated selection regime and level of recombination. Only results for surviving populations are shown. The $f_{\text{adj}}^{\text{HLA}} - f_{\text{adj}}^{\text{ALT}}$ values displayed are average values calculated over the final 2,500 y of 5,000-y simulations (Methods). SI Appendix describes the stochastic framework used for the simulations, and provides parameter definitions. We let there be five possible alleles at each of the two HLA loci. Parameter values were as follows: $b = 0.07$; $m = 0.0001$; $\varphi = 0.0015$; $\alpha = 0.002$; $Q = 5$; $\theta = 1.1$; $\Omega = 0.1$; $C = 2,000$; $k = 0.004$. r was varied between 0 and 0.06, as indicated in the x axes of each panel. For “no selection,” $\varpi = 0$ and $d = 0$; for “selection on alternative locus,” $\varpi = 0$ and $d = 0.007$; and for “pathogen selection on HLA,” $\varpi = 0.012$ and $d = 0$.

Future Directions

The system we present is necessarily a minimal caricature of the MHC, and suffers from a number of limitations. Most importantly, we have only considered the effects of interaction with a single pathogen. However, though a single HLA locus undoubtedly presents peptides from a variety of pathogens (as well as self), the selective pressure upon it will mainly arise from the pathogens causing the highest mortality. Take, for example, an HLA system as described by Fig. 1 and assume that it is under assault from n pathogens whose allelic variants may be represented according to the convention we have established as (a_i, b_i) and (x_i, y_i) ; if the most deleterious pathogen adopts the configuration (a_x, b_y) , then the homozygotes that are most disadvantaged will still be AX/AX and BY/BY.

A second important limitation of this model is that specific host recognition loci “target” specific pathogen epitopes—in other words, why should all variants at locus 1 of the pathogen specifically be recognized by locus 1 within the host? When considering associations between class I and class II HLAs, it seems justifiable to assume that different epitopes from any given pathogen are displayed by each, but it may not be strictly correct to distinguish between class I loci (particularly A and B) on this basis.

Future work in this area should also place the HLA in its wider genomic context. The very architecture of the MHC will have an effect: in the chicken, for example, the relative proximity of the TAP and class I MHC loci may have led to tight coevolution between them, limiting the possible coexpression of class I genes (20). Furthermore, in humans, HLAs interact directly with a

second family of immune system genes: Killer-cell Ig-like receptors (KIRs). KIRs display a striking haplotypic structure (21); particular KIR/HLA genotypes have been associated with different infectious disease outcomes (22, 23), and a direct effect of KIR/HLA coevolution on HLA haplotypes has recently been suggested (24).

If proven to be robust, this framework may, in principle, be able to assist in developing functional classifications of HLA alleles. It is possible to categorize HLA alleles into broad “supertypes,” based on their binding properties (25); at the same time, it is clear that a very small change in sequence (e.g., a single amino acid) can have very significant functional consequences (26). Furthermore, the ability of an HLA to bind to a specific pathogen epitope is not in itself a guarantee of an effective T-cell response to that epitope (27). If nonoverlapping allelic patterns are a signature of disease selection, they offer an alternative evolutionary approach to solving this problem. The multilocus framework described here provides a flexible platform for investigating the population-level consequences of interactions between diverse immune system genes and the pathogens they help recognize.

Methods

Deterministic Model. We used a system of linked ordinary differential equations to capture both the population genetics of the host and the disease dynamics of the pathogen. A range of coevolutionary frameworks have been developed to combine population genetics and epidemiology (28–30); the differential equation approach, first used by Gupta and Hill (31), offers a highly flexible framework that is especially amenable to the inclusion of immunological memory.

The pathogen population was represented by four potential strains ($P = 1-4$) defined by two antigenic loci containing epitopes (a, b) and (x, y) respectively (SI Appendix, Table S1). Our host population was diploid, possessing recognition alleles at two linked loci (A, B) and (X, Y), making up four possible host haplotypes ($h = 1-4$; SI Appendix, Table S2) and giving 10 possible host genotypes ($i = 1-10$; SI Appendix, Table S3). To mount an immune response against a pathogen epitope represented by a particular lowercase letter, a host must possess the recognition allele represented by the corresponding uppercase letter. The various combinations of epitopes, E_j , to which a host could be immune are shown in SI Appendix, Table S4; of these, only a subset $\{E_k\}^1$ will be accessible to host genotype i (e.g., a host genotype AXAX can only be immune to epitope sets E_1, E_3 , or E_5). A host immune to the epitopes in E_j can be infected by any pathotype not displaying those epitopes. Hosts immune to the epitopes in E_k can become immune to the epitopes of E_j by being infected by strain p , where strain p contains epitopes in E_j but not in E_k .

The dynamics of this system can be described by the following set of equations:

$$\frac{dN_i^j}{dt} = \alpha_j \omega_i + (1 - \alpha_j) \sum_{p,k} (\lambda_p N_k^i) - \left(\sum_{q,q \neq v} \lambda_q + \mu_1 \right) N_i^j - \delta_i \mu_2 G_i^j$$

$$\frac{dG_i^j}{dt} = \lambda_v (N_i^j - G_i^j) - (\sigma + \mu_1 + \mu_2) G_i^j$$

$$\frac{dI_u}{dt} = \lambda_u S_u - (\sigma + \mu_1) I_u$$

Here, N_i^j is the number of hosts of genotype i who are immune to the set of epitopes E_j . G_i^j is the number of these hosts who are infected with strain v , to which they can never mount an immune response and from which they risk dying at a rate μ_2 ; this only applies to homozygous hosts in this system (Fig. 1), so $\delta_i = 0$ for all heterozygous host genotypes. I_u is the number of hosts who are currently infected with pathogen strain u and will become immune to at least one of the epitopes of strain u . S_u is the sum of all those hosts who are not yet immune to strain u but are capable of becoming immune to at least one of the epitopes of u . All individuals recover from infection at rate σ and suffer a natural mortality rate μ_1 .

The force of infection with strain p is $\lambda_p = \frac{\beta(\lambda_p + G_p)}{\sum_{i,j} N_i^j}$, where β is a transmission coefficient, such that $R_0 = \frac{\beta}{\sigma + \mu_1}$ for the pathogen in a population of hosts that can mount an immune response against it, and $R_0 = \frac{\beta}{\sigma + \mu_1 + \mu_2}$ in a population of hosts that cannot mount an immune response against it. In the figures and figure legends, we always quote R_0 values for a pathogen in a host population that can mount an immune response against it.

Pathogen mutation can be included in the model by allowing small perturbations in the force of infection. In the model presented here we included pathogen mutation at rate m by adjusting the force of infection term, thus

$$\lambda_p^m = (1 - m)\lambda_p + \frac{1}{3} \sum_{q \neq p} \lambda_q$$

The term ω_i represents the births into the fully susceptible compartment of genotype i (thus if $j = 0$, $\alpha_j = 1$, if $j > 0$, $\alpha_j = 0$). The birth term for host genotype i is given by the following:

1. Hill AVS, et al. (1991) Common west African HLA antigens are associated with protection from severe malaria. *Nature* 352(6336):595-600.
2. Carrington M, O'Brien SJ (2003) The influence of HLA genotype on AIDS. *Annu Rev Med* 54:535-551.
3. Robinson J, et al. (2011) The IMGT/HLA database. *Nucleic Acids Res* 39(Database issue, SUPPL. 1):D1171-D1176.
4. Jeffery KJM, Bangham CRM (2000) Do infectious diseases drive MHC diversity? *Microbes Infect* 2(11):1335-1341.
5. Hedrick PW (2002) Pathogen resistance and genetic variation at MHC loci. *Evolution* 56(10):1902-1908.
6. Trowsdale J (2011) The MHC, disease and selection. *Immunol Lett* 137(1-2):1-8.
7. Cao K, et al. (2001) Analysis of the frequencies of HLA-A, B, and C alleles and haplotypes in the five major ethnic groups of the United States reveals high levels of diversity in these loci and contrasting distribution patterns in these populations. *Hum Immunol* 62(9):1009-1030.
8. Cao K, et al. (2004) Differentiation between African populations is evidenced by the diversity of alleles and haplotypes of HLA class I loci. *Tissue Antigens* 63(4):293-325.
9. Shaw CK, Chen LL, Lee A, Lee TD (1999) Distribution of HLA gene and haplotype frequencies in Taiwan: A comparative study among Min-nan, Hakka, Aborigines and Mainland Chinese. *Tissue Antigens* 53(1):51-64.

$$\omega_i = \kappa f_h f_g (1 + \delta_i),$$

where κ is the total death rate for the entire population; δ_i is defined as above, and f_h and f_g are the frequencies of the haplotypes that make up host genotype i .

Haplotype frequencies are calculated as follows, where r is the host recombination rate. If $r = 0.5$, the two host loci are effectively unlinked.

$$f_h = \frac{2 \sum_j N_j^{c1} + \sum_j N_j^{c2} + \sum_j N_j^{c3} + (1-r) \sum_j N_j^{c4} + r \sum_j N_j^{c5}}{2 \sum_{i,j} N_i^j}$$

See SI Appendix, Table S5 for the values of c_{1-5} that correspond to a particular haplotype.

The total death rate is calculated as follows:

$$\kappa = \mu_1 \left(\sum_{i,j} N_i^j \right) + \sum_{i,j} G_i^j \mu_2$$

Numerical simulations were carried out using the ode45 solver in MatLab version 7.10.0 (R2010b).

Stochastic Model. A full description of the stochastic model is provided in SI Appendix, section 1. Briefly, the population was made up of N hosts, where $N < C$, the population carrying capacity. Each host was represented by a 19-element identifier code that recorded age, genotype, infection, and immunity status. As in the deterministic model, host genotype AX/AX was only capable of becoming immune to pathogen epitopes a and x , and risked death when infected with a pathogen it could not recognize. Infection, recovery, mortality, and reproduction were all probabilistic events.

Metrics. We used a standard metric (Lewontin's D' , normalized where necessary for >2 alleles per locus, as described in ref. 17) to measure LD.

The f^* metric for nonoverlap between two loci was calculated as described in ref. 18 and adjusted as follows:

$$f_{adj}^* = (1 - H_{max}) f^*,$$

where H_{max} is the frequency of the most frequent haplotype in the population. f^* takes values between 0 and 1, where values closer to 1 indicate a more nonoverlapping pattern. However, $f^* = 1$ for a population that consists of one haplotype only, which is not a case of true nonoverlap. For f_{adj}^* , by contrast, populations containing relatively balanced frequencies of non-overlapping haplotypes will receive the highest scores.

To calculate $f_{adj}^* HLA - f_{adj}^* ALT$ from our simulations in Fig. 5, we measured $f_{adj}^* HLA - f_{adj}^* ALT$ every 20 y during the final 2,500 y of a 5,000-y simulation, and took the mean of those measurements.

ACKNOWLEDGMENTS. We thank Adrian Hill, Angus Buckling, Paul Harvey, and Oliver Pybus for their comments on the manuscript, and Adrian Smith for general guidance on this project. Funding for this work was provided by the Wellcome Trust, the European Research Council (ERC Advanced Grant - DIVERSITY), the Biotechnology and Biological Sciences Research Council, and the Christopher Welch Trust. B.S.P. is a Sir Henry Wellcome Postdoctoral Fellow (Grant 096063/Z/11/Z) and a Junior Research Fellow at Merton College, Oxford. S.G. is a Royal Society Wolfson Research Fellow and an ERC Advanced Investigator.

10. Cox ST, et al. (1999) HLA-A, -B, -C polymorphism in a UK Ashkenazi Jewish potential bone marrow donor population. *Tissue Antigens* 53(1):41-50.
11. Buhler S, Nunes JM, Nicoloso G, Tiercy JM, Sanchez-Mazas A (2012) The heterogeneous HLA genetic makeup of the Swiss population. *PLoS ONE* 7(7):e41400.
12. Mohyuddin A, et al. (2002) HLA polymorphism in six ethnic groups from Pakistan. *Tissue Antigens* 59(6):492-501.
13. Gupta S, Ferguson N, Anderson R (1998) Chaos, persistence, and evolution of strain structure in antigenically diverse infectious agents. *Science* 280(5365):912-915.
14. Anderson RM, May RM (1991) *Infectious Diseases of Humans: Dynamics and Control* (Oxford Univ Press, New York).
15. Kouyos RD, Salathé M, Otto SP, Bonhoeffer S (2009) The role of epistasis on the evolution of recombination in host-parasite coevolution. *Theor Popul Biol* 75(1):1-13.
16. Carrington M (1999) Recombination within the human MHC. *Immunol Rev* 167:245-256.
17. Hedrick PW (1987) Gametic disequilibrium measures: Proceed with caution. *Genetics* 117(2):331-341.
18. Buckee CO, Gupta S, Kriz P, Maiden MCJ, Jolley KA (2010) Long-term evolution of antigen repertoires among carried Meningococci. *Proc R Soc B Biol Sci* 277(1688):1635-1641.
19. Weitkamp LR, Ober C (1999) Ancestral and recombinant 16-locus HLA haplotypes in the Hutterites. *Immunogenetics* 49(6):491-497.

Supporting Information Appendix

Pathogen selection drives nonoverlapping associations between HLA loci.

Bridget S. Penman^a, Ben Ashby^a, Caroline O. Buckee^b and Sunetra Gupta^{a,1}

Author affiliations:

- a. Department of Zoology, University of Oxford, South Parks Road, Oxford, OX13PS, UK
- b. Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard School of Public Health, Boston, USA, 02115.

Corresponding author:

1. Sunetra Gupta, Department of Zoology, University of Oxford, South Parks Road, Oxford, OX13PS UK
sunetra.gupta@zoo.ox.ac.uk

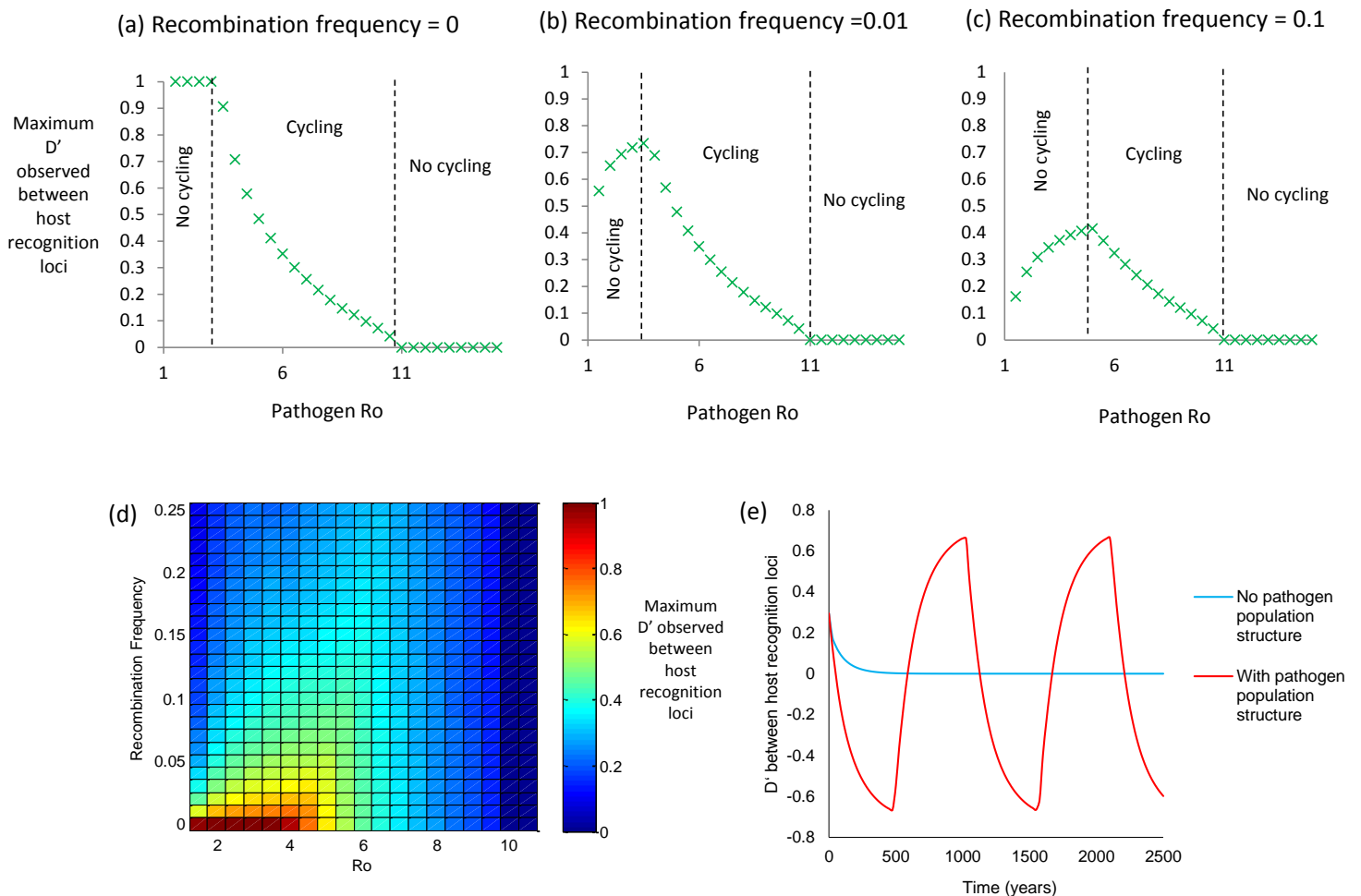
This appendix contains the following:

| | |
|---|---|
| 1. A sensitivity analysis of the deterministic HLA model | 2 |
| 2. The stochastic HLA model and its behaviour | 3 |
| 2.1 The stochastic model..... | 3 |
| 2.2 Equivalence between the stochastic and deterministic HLA models | 5 |
| 2.3 Extension to a 3 locus, 3 allele system..... | 6 |
| 3. Supplementary Tables | 7 |

1. A sensitivity analysis of the deterministic HLA model

As noted in the main text, cyclical and stable structuring of both pathogen and host emerge within our framework as a result of immunological feedbacks. The nature of this structuring, however, is sensitive to the basic reproductive number of the pathogen (R_0) and the recombination frequency between host recognition loci. Figure S1 explores these effects, using a standard measure of linkage disequilibrium (Lewontin's D' [17]) to measure the degree of structuring that is generated between the host loci.

Figure S1: The effects of varying R_0 and recombination frequency on host genetic structuring. Panels (a-c) illustrate the effects of varying R_0 and the recombination frequency on the maximum D' observed between host loci. In these panels, host mortality rate (μ_1)=0.05 ; mutation rate (m) =0.00001; pathogen mortality rate (μ_2)=5 and recovery rate (σ) =10. The heatmap in panel (d) illustrates the maximum D' possible for different levels of R_0 combined with different recombination frequencies, when host mortality rate (μ_1)=0.05; mutation rate (m) =0.00001; pathogen mortality rate (μ_2)=5 and host recovery rate (σ) =7. The time series in panel (e) compares the change in D' over time for a host where the pathogen population is allowed to become structured (as described in the main text), with a host whose pathogen population is forced to remain unstructured (i.e. all possible pathogen strains are always present). For this panel, host mortality rate (μ_1)=0.05; mutation rate (m) =0.0001; pathogen mortality rate (μ_2)=4; host recovery rate (σ) =8; recombination frequency (r)=0.01 and R_0 =4.75.

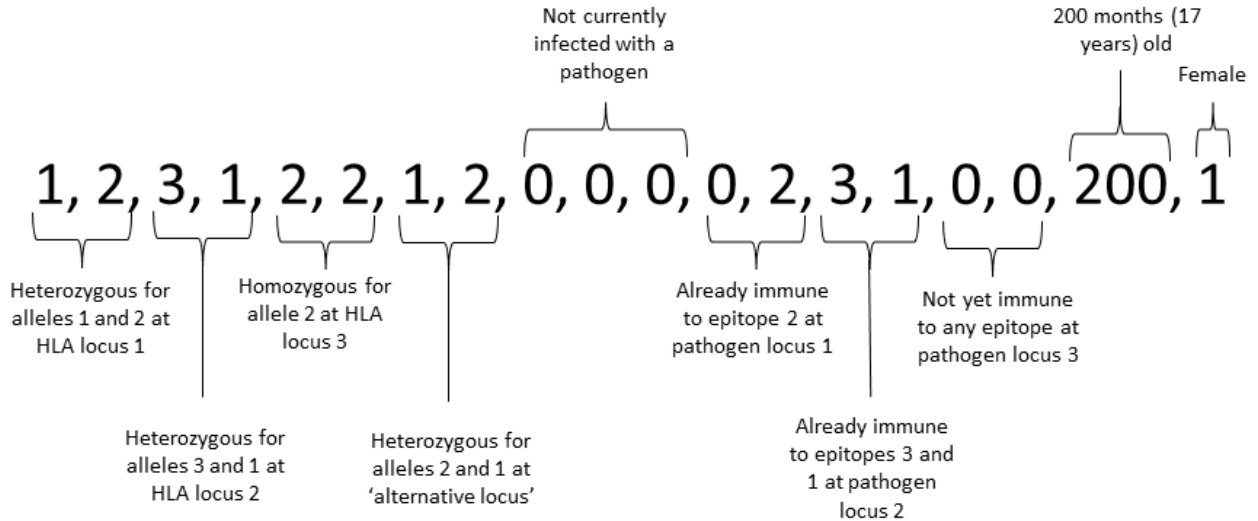


2. The stochastic HLA model and its behaviour

2.1 The stochastic model

We represent each member of the host population by a 19 element identifier code (H_1 - H_{19}), recording: (i) their genotype at up to 3 different HLA loci (H_1 - H_6); (ii) their genotype at an alternative locus, linked to the HLAs, that can be under non-HLA forms of selection (H_7 - H_8) ; (iii) the genotype of up to one pathogen infecting them (H_9 - H_{11}); (iv) whether or not they have generated a memory immune response against pathogen epitopes they have already been exposed to (H_{12} - H_{17}); (v) their age (H_{18}) and finally (vi) their gender (H_{19}).

A typical host identifier might appear as follows:



At each of the loci making up the host genotype, the first element represents the maternally derived allele and the second element represents the paternally derived allele. We are thus able to keep track of haplotypes.

As in the deterministic model described in the main text, the variants found at a particular HLA locus specifically target the variants of a particular pathogen epitope locus. This stochastic framework allows us to define any number of variants at each pathogen locus, and equivalent recognition alleles in the host. Section 2.3 of this document illustrates the 3 allele, 3 locus case.

At every time step of the simulation, the following events occur:

- Individuals age by one month.
- Uninfected individuals become infected with strain $ijkl$, with probability p_{jkl}

$$\text{if } H_{12,13} \neq j, H_{14,15} \neq k \text{ and } H_{16,17} \neq l \quad p_{jkl} = \frac{\sum_{x\dots z} H_{x\dots z}}{\sum_{H_7=j; H_8=k; H_9=l} H_{x\dots z}} \cdot \theta$$

$$\text{otherwise} \quad p_{jkl} = 0$$

where $H_{x\dots z}$ = any host within the population and θ controls the likelihood of transmission per infection. θ is directly equivalent to the transmission parameter β in the deterministic model.

- For every infection that occurs, mutation of a randomly chosen pathogen epitope (to any of the alleles that are possible at that locus) occurs with probability m .
- Infected individuals recover with probability Ω . The average length of infection within this system is therefore $\frac{1}{\Omega}$ months.
- Upon recovery from infection with $ijkl$, if $H_1 = j$ then $H_{12} = j$; if $H_2 = j$ then $H_{13} = j$; if $H_3 = k$ then $H_{14} = k$; if $H_4 = k$ then $H_{15} = k$; if $H_5 = l$ then $H_{16} = l$; if $H_6 = l$ then $H_{17} = l$.
- Infected individuals who can recognise none of the pathogen's epitopes (i.e. for whom $H_{10,11} \neq j$, $H_{12,13} \neq k$ and $H_{14,15} \neq l$) die with probability ϖ . Upon death, individuals are removed from the population.
- If we are imposing selection at the 'alternative locus' as a substitute for HLA interacting pathogen selection (see figure 5 in the main text), a maximum of 2 out of 4 possible alleles that may occupy the alternative locus are designated 'favoured' at any one time. Every time step, there is a probability k that the identity of an allele occupying one of the 2 'favoured' slots will change. If a change occurs, there is an equal chance (1/5) that the newly chosen 'favoured' allele will be any one of the alleles that can exist at the locus, or 'null' – i.e. no favoured allele at that time. Any individual who lacks a favoured allele at their alternative locus has a probability d of dying within a given time step.
- Females over the age of 15 years (180 months) reproduce with a randomly chosen male partner (also over the age of 15 years), with probability b . A new host (aged 0 months, with a genotype generated from a combination of maternal and paternal HLA haplotypes) is added to the population. Intra-haplotypic recombination is assumed to occur between the HLA loci, independently in both parents, with probability r . If the population is already at carrying capacity (C), the new host displaces a randomly chosen existing host.
- Q new individuals, with randomly generated genotypes, are introduced to the population with probability α ; replacing Q existing individuals in the population. This step simulates gene flow between the simulated population and the wider world.
- Random host death occurs with probability ϕ .
- All hosts over the age of 40 years (480 months) are removed from the population.

A programme to perform these operations was written in C. A Mex file was created so that the model could be called from within Matlab, version 7.10.0 (R2010b).

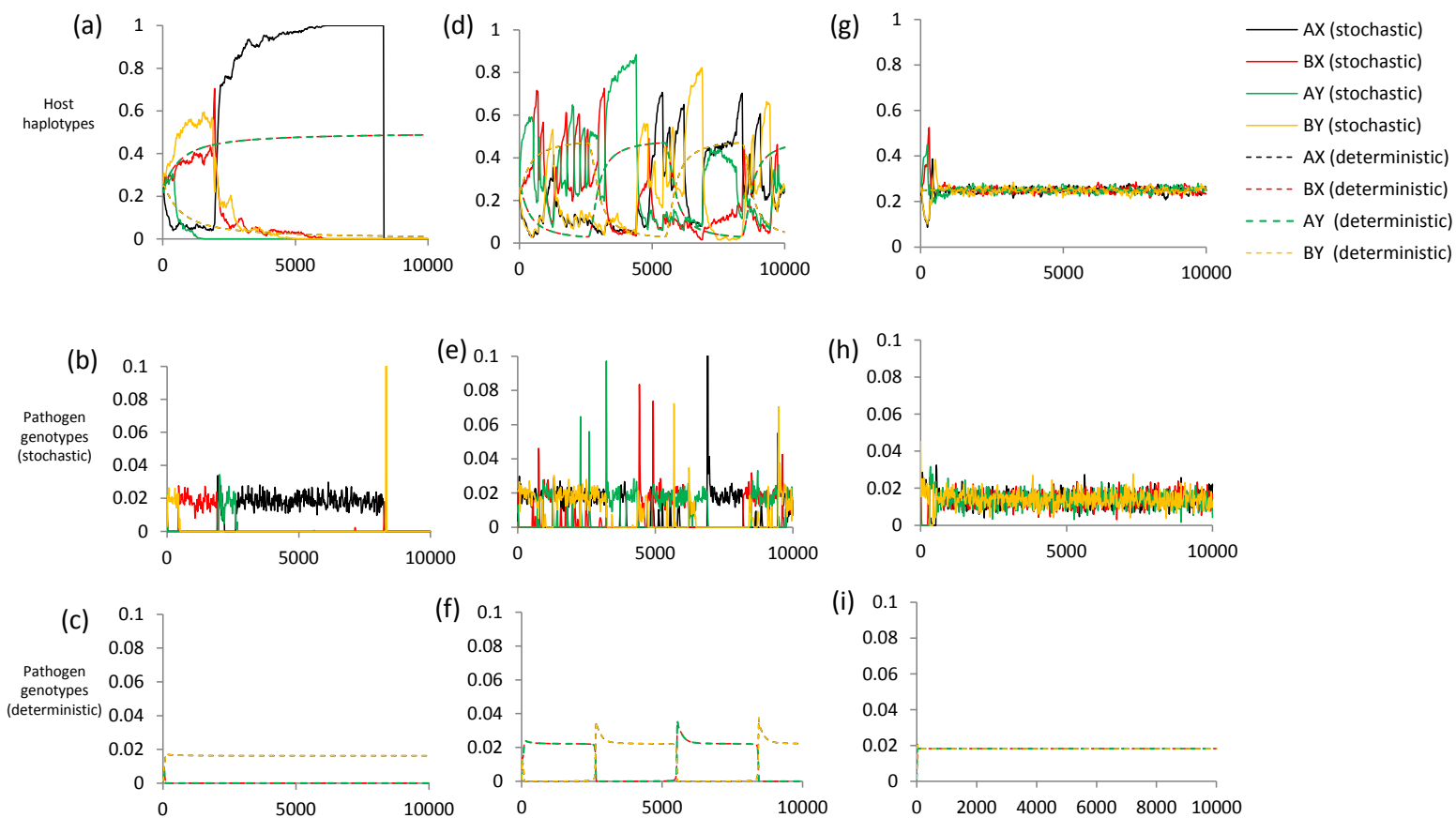
The initial conditions for each of our stochastic simulations were either (i) a population containing entirely randomly generated host and pathogen genotypes, or (ii) a population in which a certain proportion of host haplotypes were a specific 'founder' haplotype, and the rest were randomly generated. In these latter simulations, the *initial* pathogen population was such that its epitopes could definitely be recognised by the 'founder' haplotype. This condition was intended to avoid the host population being wiped out by a pathogen within a very short space of time, but the pathogen could quickly mutate to adopt alternative

structures. For figure 5 of the main text, the initial frequency of the founder haplotype was 0.505 in all simulations.

2.2 Equivalence between the stochastic and deterministic HLA models

When we consider a 2 locus, 2 allele recognition system and only apply HLA-interacting pathogen selection (i.e. $d=0$; $\varpi>0$), the stochastic model behaves equivalently to the deterministic model described in the main text.

Figure S2: Comparing the stochastic and deterministic models, in a 2 locus, 2 allele system. The behaviour of the system is illustrated under 3 different sets of conditions: those which lead to permanent genetic structuring in the deterministic model (panels a-c); those which lead to cyclical behaviour in the deterministic model (panels d-f), and those which lead to no structuring in the deterministic model (panels g-i). For the deterministic model, $\mu_1 = 0.04$ years⁻¹; $\mu_2 = 0.5$ years⁻¹; $r = 0$; $m = 0.0001$; $\sigma = 1.2$ years⁻¹; $\beta = 2.5$ in (a-c); 4 in (d-f) and 10 in (g-i). For the stochastic simulations, $\phi = 0.0035$; $b = 0.03$; $Q = 0$; $\Omega = 0.1$; $\varpi = 0.04$; $m = 0.0001$; $r = 0$; $C = 1500$; $\alpha = 0$; $d = 0$; $k = 0$; $\theta = 2.5$ in (a-c); 4 in (d-f) and 10 in (g-i). In all panels, the simulated populations contained randomly generated host and pathogen genotypes at time = 0.

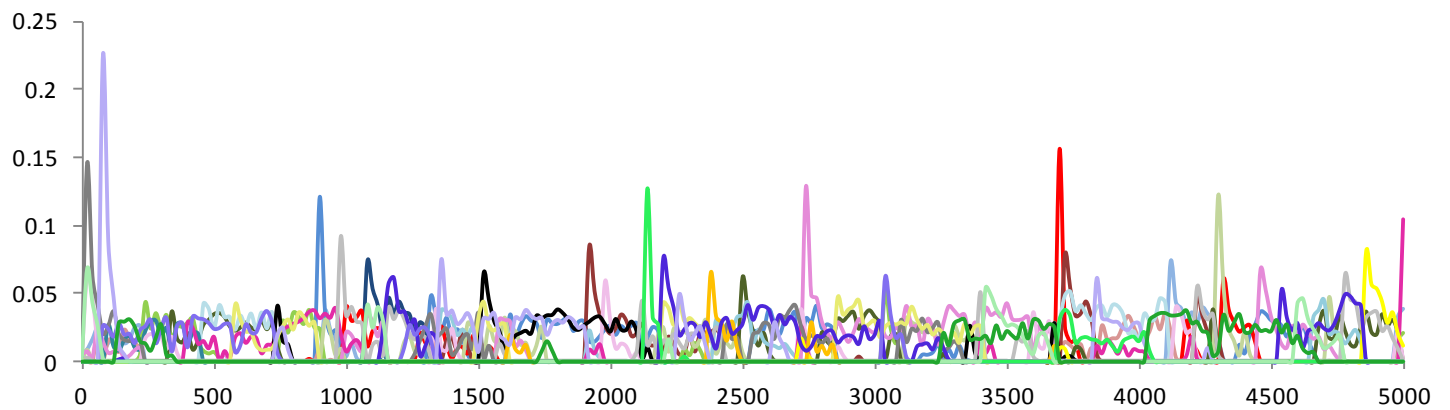


2.3 Extension to a 3 locus, 3 allele system

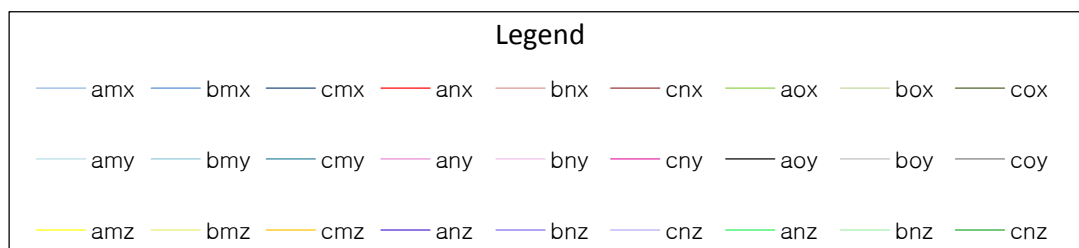
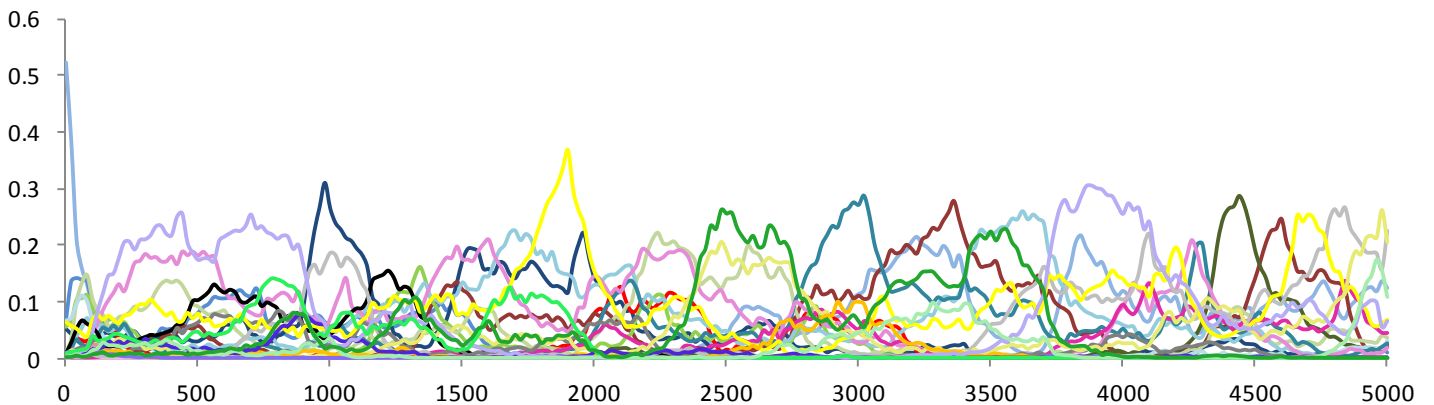
The stochastic framework allows us to extend the principles outlined in figure 1 of the main text to a system of 3 pathogen epitope loci with 3 alleles at each (a,b,c ; m,n,o and x,y,z), and 3 corresponding host recognition loci (A, B, C; M, N, O and X, Y, Z). The cyclical behaviour reported in the 2 locus, 2 allele model can still be observed in this higher-dimension system (figure S3).

Figure S3: Pathogen selection and host haplotypes structuring in a 3 locus, 3 allele stochastic system. Panel (a) illustrates the proportion of the population infected with any one of 27 pathogen strains made up of 3 epitopic loci, and panel (b) illustrates the behaviour of host haplotypes in the same population. Equivalent colours have been used for pathogen genotypes and host haplotypes, thus the colour indicated for pathogen 'anx' also refers to host haplotype 'ANX'. In all panels, $b=0.05$; $m=0.0003$; $r=0.005$; $\phi=0.0035$; $\Omega=0.1$; $\alpha=0.0001$ $C=1500$; $Q=5$; $\theta=1.5$; $\varpi=0.01$. The simulated population began with a 'founder' host haplotype of 'AMX' at a frequency of 0.52, whilst all other host haplotypes at time=0 were randomly generated.

(a) Proportion infected



(b) Host haplotype frequencies



3. Supplementary Tables

Table S1: Pathogen strains

| Pathotype subscript (p) | Epitopes |
|-----------------------------|----------|
| 1 | ax |
| 2 | bx |
| 3 | ay |
| 4 | by |

Table S2: Host haplotypes

| Haplotype subscript (h) | Haplotype |
|-----------------------------|-----------|
| 1 | AX |
| 2 | BX |
| 3 | AY |
| 4 | BY |

Table S3: Host genotypes

| Genotype superscript (i) | Genotype |
|------------------------------|----------|
| 1 | AXAX |
| 2 | BXBX |
| 3 | AYAY |
| 4 | BYBY |
| 5 | AXBY |
| 6 | BXAY |
| 7 | AXBX |
| 8 | AXAY |
| 9 | BXBY |
| 10 | AYBY |

Table S4: Possible sets of epitopes a host could be immune to

| j | E_j |
|-----|------------|
| 0 | None |
| 1 | a |
| 2 | b |
| 3 | x |
| 4 | y |
| 5 | a, x |
| 6 | b, x |
| 7 | a, y |
| 8 | b, y |
| 9 | a, x, y |
| 10 | b, x, y |
| 11 | a, b, x |
| 12 | a, b, y |
| 13 | a, b, x, y |
| 14 | a, b |
| 15 | x, y |

Table S5 Host genotype identifiers (c_{1-5}) for calculating the frequency of each haplotype

| Haplotype subscript (h) | c_1 | c_2 | c_3 | c_4 | c_5 |
|-----------------------------|-------|-------|-------|-------|-------|
| 1 | 1 | 7 | 8 | 5 | 6 |
| 2 | 2 | 7 | 9 | 6 | 5 |
| 3 | 3 | 8 | 10 | 6 | 5 |
| 4 | 4 | 9 | 10 | 5 | 6 |