

## Neutral genomic signatures of host-parasite coevolution

Daniel Živković<sup>a,\*</sup>, Sona John<sup>a</sup>, Mélissa Verin<sup>a,b</sup>, Wolfgang Stephan<sup>c</sup>, Aurélien Tellier<sup>a,\*</sup>

<sup>a</sup>*Section of Population Genetics, Technical University of Munich, 85354 Freising, Germany*

<sup>b</sup>*Department of Mathematics and Statistics, Queen's University, Ontario K7L 3N6, Canada*

<sup>c</sup>*Leibniz Institute for Evolution and Biodiversity Science, 10115 Berlin, Germany*

---

\*Corresponding authors.

*Email addresses:* daniel.zivkovic@gmx.at, tellier@wzw.tum.de

## 1 **Abstract**

2 Coevolution is a selective process of reciprocal adaptation between antagonistic or mutualistic sym-  
3 bionts and their host. Classic population genetics theory predicts the signatures of selection at the  
4 interacting loci but not the neutral genome-wide polymorphism patterns. We here build a coevo-  
5 lutionary model with cyclic changes in the host and parasite population sizes. Using an analytical  
6 framework, we investigate if and when these population size changes can be observed in the neu-  
7 tral site frequency spectrum of the host and parasite full genome data. We show that polymorphism  
8 data sampled over time can capture the changes in the population size of the parasite but not of the  
9 host because genetic drift and mutations occur on different time scales in the coevolving species.  
10 This is due to the small parasite population size at the onset of the coevolutionary history subse-  
11 quently undergoing a series of strong bottlenecks. We also show that tracking coevolutionary cycles  
12 is more likely for a small amount of parasite per host and for multiple parasite generations per host  
13 generation. Our results demonstrate that time sampling of host and parasite full genome data are  
14 crucial to infer the co-demographic history of interacting species.

## 15 **Keywords**

16 population genomics; epidemiological model; host-parasite coevolution; SI model; population  
17 demographic history; population dynamics

## Introduction

1

2 Host-parasite antagonistic interactions are a role model for observing and studying rapid evolu-  
3 tionary change as well as feedbacks between ecological and evolutionary forces and time scales.  
4 Coevolution, defined here as the reciprocal adaptation of hosts and their parasites, typically gener-  
5 ates significant phenotypic and genetic diversity for host resistance and for parasite infectivity and  
6 virulence. Such changes in the genetic composition of the interacting species at the key underpin-  
7 ning loci, drive, and are driven by, short-term epidemiological (ecological) dynamics. To develop  
8 infectious disease epidemiology as a predictive science, there is thus a need to understand the  
9 synergy of fast evolution and within and between populations disease dynamics [1], the so-called  
10 eco-evolutionary feedbacks [2].

11 Coevolution as determined by changes in allele frequencies over time at the interacting genes, is  
12 observable as coevolutionary cycles driven by negative indirect frequency-dependent selection [3,  
13 4]. Theory predicts that a continuum of dynamics of allele frequency cycles occurs, characterized  
14 by their stability, period and amplitude, and ranging between two extremes: the arms race and the  
15 trench warfare dynamics. The arms race is defined as the recurrent fixation of alleles at these major  
16 loci in host and parasite populations [5, 6], while trench warfare maintains cycling over a long  
17 period of time [7] (also called the Red Queen dynamics [6], or Fluctuation Selection Dynamics  
18 [2]). The transition between these types of dynamics depends on the occurrence and strength  
19 of negative direct frequency-dependent selection [4], which stabilizes cycles and is generated by  
20 several host and parasite life history traits (reviewed in [8]).

21 Theory also predicts that the coevolutionary dynamics, either by arms race or trench warfare,  
22 can be observed in the patterns of polymorphism at these loci, namely in the frequencies of Sin-  
23 gle Nucleotide Polymorphisms (SNPs). The arms race is expected to generate recurrent selective  
24 sweeps, while trench warfare generates balancing selection. These predictions form the basis of  
25 scans for genes under coevolution in host or parasite genomes relying on the prevalent perception  
26 that natural selection acts only at few loci, while neutral forces, such as demographic histories,  
27 affect the whole genome. Detecting genes under coevolution entails therefore to disentangle the  
28 signatures of arms race or trench warfare from the polymorphism patterns observed in genome-  
29 wide data.

30 Besides the allelic coevolutionary cycles at the host and parasite interacting loci, size fluctua-  
31 tions of host and parasite populations are also predicted to occur. These changes in population

1 size over time are induced by reciprocal selection among the antagonists and are an inherent prop-  
2 erty of host-parasite coevolution under epidemiological (or Lotka-Volterra) dynamics (such as the  
3 Susceptible-Infected or Susceptible-Infected-Recovered models, [9]). In a more complex coevolu-  
4 tionary system with several host and parasite genotypes being present at major genes of interac-  
5 tion, cycles of coevolution do occur, thereby generating a fluctuation of the numbers of hosts and  
6 parasites over time [9] as an epidemiological feedback [2]. Several episodes of coevolution pro-  
7 ceed with increasing and decreasing disease prevalence depending on the cycling of resistance and  
8 infectivity alleles. The epidemiological feedback generates negative direct frequency-dependent  
9 selection, thus stabilizing the frequencies of alleles and maintaining long-term diversity at the  
10 interacting loci [10]. Coevolutionary models based on Lotka-Volterra dynamics have similar char-  
11 acteristics [11–14]. Population size changes due to coevolution should affect the whole genome  
12 polymorphism of both antagonistic species, an effect which we term as the co-demographic history.  
13 When studying host and parasite polymorphism data, two sources of demographic variation gener-  
14 ating genetic drift can therefore be defined: 1) the population or species demographic history  
15 (e.g. colonisation of new habitats or recolonisation), and 2) the co-demographic history due to co-  
16 evolutionary and epidemiological dynamics. Both types of demographic events affect the ability to  
17 detect genes under coevolution using scans for arms race or trench warfare signatures. Moreover,  
18 there is currently no theoretical prediction regarding the signature of co-demographic history on  
19 genome-wide polymorphism in hosts and parasites.

20 Our aim in this study is to propose the first model of neutral polymorphism generated by the co-  
21 demographic history of host and parasite populations. First, we establish an epidemiological model  
22 describing changes in the numbers of healthy and infected hosts over time focusing on biallelic  
23 gene-for-gene and matching-allele infections and initially assigning one parasite per host. Second,  
24 we utilize an analytical result [15] for the neutral site frequency spectrum (SFS) under arbitrary  
25 deterministic population size changes and apply it to the host and parasite populations. We show  
26 that these population size changes can be quite drastic in the parasite and occur on a time scale  
27 slow enough to leave a corresponding signature in the SFS over time. Conversely, changes in the  
28 host size are barely detected in the polymorphism data. Finally, since such recurrent bottlenecks  
29 in parasites cause a reduced amount of polymorphism, we further discuss the impacts of multiple  
30 parasites per host and multiple parasite generations per host generation.

# Theoretical and computational framework

## Modeling a single coevolving locus

We have a haploid one locus model with  $A$  alleles in the host and in the parasite. The infection matrix is given by  $\mathcal{A} = (\alpha_{ij})$  with  $1 \leq i, j \leq A$ . The entries  $\alpha_{ij}$  give the probability that once encountered hosts of genotype  $i$  are infected by parasites of genotype  $j$ . Examples of simple infection matrices for two alleles are given in Table 1.

In analogy to [16] the changes of host and parasite allele frequencies over time are determined by the following coupled differential equations:

$$\frac{dH_i}{dt} = H_i \left[ b_i(1 - c_{H_i}) - d_i - \sum_{j=1}^A \alpha_{ij} \beta_{ij} (1 - c_{P_j}) \sum_{k=1}^A I_{kj} \right] + b_i(1 - c_{H_i}) \sum_{j=1}^A (1 - s_{ij}) I_{ij}, \quad (1)$$

$$\frac{dI_{ij}}{dt} = I_{ij}(-d_i - \delta_{ij}) + H_i \left[ \alpha_{ij} \beta_{ij} (1 - c_{P_j}) \sum_{k=1}^A I_{kj} \right]. \quad (2)$$

In Equations (1) and (2),  $H_i$  is the number of healthy (*i.e.* non-infected) individuals of genotype  $i$ , and  $I_{ij}$  denotes the number of host genotype  $i$  infected by parasite genotype  $j$ .  $b_i$  and  $d_i$  are the birth and natural death rates (*i.e.* independent of the disease) of host genotype  $i$ , respectively, and  $\delta_{ij}$  is the disease induced death rate caused by pathogen  $j$  on host genotype  $i$  (*i.e.* the effect of pathogen on host mortality [16]).  $\beta_{ij}$  is the disease transmission rate between a parasite of genotype  $j$  and a host of genotype  $i$ .  $c_{H_i}$  and  $c_{P_j}$  are the costs for the hosts and the parasites of carrying genotype  $i$  and  $j$ , respectively.  $s_{ij}$  is the decrease of reproductive fitness of host genotype  $i$  due to an infection of parasite  $j$ , *i.e.* the effect of pathogen on host fecundity. Due to the large number of parameters

Table 1 Infection matrices for four coevolution models

matching-allele	inverse matching-allele	gene-for-gene	inverse gene-for-gene
$\begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$	$\begin{pmatrix} 0 & 1 \\ 1 & 1 \end{pmatrix}$	$\begin{pmatrix} 1 & 0 \\ 0 & 0 \end{pmatrix}$

The infection matrices determine the outcome of the interaction between host genotypes (rows) and parasite genotypes (columns). To keep the illustration simple, the rates  $\alpha_{ij}$  are either chosen as one for infection or as zero for resistance. Matching-allele and gene-for-gene models are shown as well as their inverse versions.

1 in our epidemiological model, we investigate a simplified version with two alleles ( $A = 2$ ) in the  
2 host and in the parasite except for the evaluation of the host effective population size over time (SI  
3 1) and the reproduction ratios (SI 2), which can be computed for arbitrary  $A$ . Due to the purely  
4 deterministic setting, an allele is considered as lost as soon as its count takes a value below one.  
5 We only allow rates and costs to differ among the genotypes in SI 1, SI 2, and for the calculation  
6 of the fixed points (SI 3) of the dynamical system (1) and (2). Otherwise, we assume  $b_i = b$ ,  
7  $d_i = d$ ,  $\delta_{ij} = \delta$ ,  $\beta_{ij} = \beta$  and  $s_{ij} = s$ . The entries  $\alpha_{ij}$  of the infection matrix are either zero or one  
8 depending on the considered model (see Table 1). While we evaluate the fixed points (SI 3) for  
9 all four coevolutionary models, only the matching-allele (MA) and the gene-for-gene (GFG) model  
10 are investigated in further detail, since the inverse matching-allele (iMA) model is symmetric (and  
11 therefore equivalent) to the MA model for  $A = 2$  and the inverse gene-for-gene (iGFG) model is  
12 restricted in its behavior compared to the GFG model.

13 For the MA model symmetric costs are chosen and (except for SI 1-SI 3) assumed to be inde-  
14 pendent of the host ( $c_{H_i} = c_H$ ) and the parasite ( $c_{P_j} = c_P$ ) genotype. For the GFG model costs are  
15 chosen asymmetrically as follows (except for the choice of arbitrary costs in SI 1-SI 3). Since  $H_1$   
16 defends itself successfully against  $P_1$ , and  $P_2$  can infect both host genotypes, only  $H_1$  and  $P_2$  are  
17 assumed to carry realistically small costs  $c_{H_1}$  and  $c_{P_2}$  ( $c_{H_2} = c_{P_1} = 0$ , see [4]).

18 The total number of hosts of genotype  $i$  is given by  $W_i = H_i + \sum_{j=1}^A I_{ij}$ . The number  $P_j$  of  
19 parasites of genotype  $j$  is only implicitly given in (1) and (2) by the number of infected individuals;  
20 assuming one parasite per infected genotype we have  $P_j = \sum_{i=1}^A I_{ij}$ .

21 The change in the effective population size of the host over time  $N^W = \sum_{i=1}^A W_i$  is obtained by  
22 numerically solving (1) and (2). The respective differential equation and a condition for obtaining  
23 a constant population size are given in SI 1. The effective population size of the parasite is obtained  
24 as  $N^P = \sum_{j=1}^A P_j$  (as we assume first one parasite per host).

## 25 **Assessing the effect of the coevolving locus on neutral polymorphism patterns**

26 To evaluate the impact of host-parasite interactions on neutrally evolving and genome-wide dis-  
27 tributed SNPs over time for interesting cases (as determined by the stability analysis), we utilize  
28 the SFS (SI 4), which is the distribution  $f_{n,i}(t)$  of the number of times  $i$  a mutant allele is observed  
29 across sites in a sample of  $n$  DNA sequences at time  $t$ . While  $f_{n,i}(t)$  takes allele counts in absolute  
30 numbers, its relative version  $r_{n,i}(t) = f_{n,i}(t) / \sum_{k=1}^{n-1} f_{n,k}(t)$  is normalized by the total number of

1 segregating sites. For this purpose, the deterministic trajectories of time-varying host and parasite  
2 population sizes are employed into the analytical result [15] for the neutral SFS. This application  
3 requires an appropriate scaling: we define a relative population size function  $\rho(t)$  as the ratio of  
4 the population size  $N(t)$  at time  $t$  scaled by the reference population size  $N_{\text{ref}}$  at the time of the in-  
5 fection, which initiates the coevolutionary history, *i.e.*  $\rho(t) = N(t)/N_{\text{ref}}$ , and denote the population  
6 mutation rate as  $\theta = 2 N_{\text{ref}} \mu$ . We compute the changes in frequencies of neutral alleles generated  
7 by the co-demographic scenario in terms of the SFS (also in relative numbers) and the average  
8 number of pairwise differences  $\Pi_n(t) = 1/\binom{n}{2} \sum_{k=1}^{n-1} k(n-k) f_{n,k}(t)$  (SI 4) for  $n = 20$  hosts and  
9 parasites. Our forward approach is adequately suited for the analysis of time-series data in contrast  
10 to the corresponding coalescent result for the neutral allelic spectrum [17, 18], where only a single  
11 time point can be immediately evaluated for a given demographic history.

## 12 **Key characteristics of our model**

13 Our model presents three key features to keep in mind when reading the results below. The first  
14 key aspect of our eco-evolutionary framework is that changes in the population size are a direct  
15 consequence of the dynamics of the model, driven by a single locus underpinning coevolution, and  
16 not assumed to follow an arbitrarily chosen function of time as in the majority of the population  
17 genetics literature.

18 The second crucial point is the definition of the host and parasite time scales of evolution as  
19 determined by the generation times and the population mutation rates of the antagonistic species  
20 [19]. If viral, bacterial or fungal parasites often have higher mutation rates than their hosts, their  
21 effective population size may not always be larger than that of the hosts and at the onset of an  
22 epidemics. The reference population size  $N_{\text{ref}}$  at the onset of an epidemics is important because 1)  
23 it sets up the initial available diversity, and 2) it defines the time scale for genetic drift in host and  
24 parasite and the timing of new neutral mutations occurring with rate  $\theta$ . In our population genetics  
25 setting, time is scaled in units of  $N_{\text{ref}}$  generations, whereas the host-parasite model specified in  
26 Equations (1) and (2) runs on arbitrary continuous time reflecting calendar time (in weeks, months  
27 or years). If calendar time is equivalent for both species, the scaled time based on  $N_{\text{ref}}$  defines the  
28 changes occurring in the observed polymorphism over time. We exemplify the difference in time  
29 scale and its influence on polymorphism data by a simplified bottleneck model. Two populations  
30 with different initial population sizes  $N_{\text{ref}}$  experience a size change of the same magnitude and for

1 the same number of generations on calendar time scale (bottom x-axes) but for different rescaled  
 2 time with respect to  $N_{\text{ref}}$  (top x-axes of Figure 1, a and c). Consequently, a size change of the  
 3 same magnitude but based on two different initial population sizes  $N_{\text{ref}}$  affects the SFS similarly  
 4 regarding the course of time but with very different strength and detectability (y-axes of Figure 1,  
 5 b and d). Note that the number of singletons scaled by  $\theta$  is used as a representation of the SFS  
 6 for illustrative purposes in this and several subsequent examples because these show the most  
 7 pronounced signals of all allelic classes, but results are similar using  $\Pi_n(t)$ .

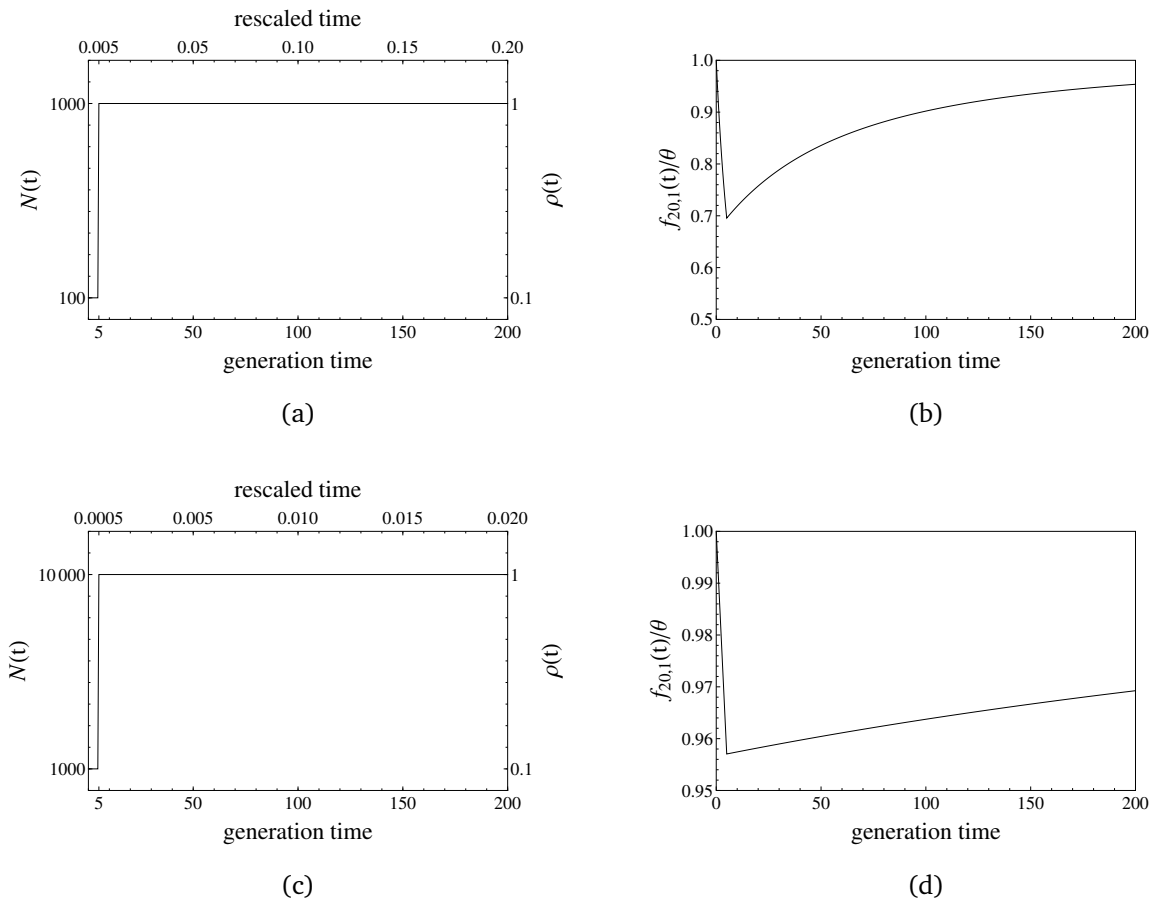


Figure 1 A simple bottleneck model with an initial population size of 1000 (a) and 10000 (c) is considered, where the population size drops instantaneously to one tenth of its original size at time zero for 5 generations before recovering again to its original size. Original time in generations is opposed to time scaled in units of the initial population size on the x-axes and the population size changes are given in absolute and relative numbers on the y-axes. The absolute number of singletons scaled by  $\theta$ ,  $f_{20,1}(t)/\theta$ , are illustrated for the respective cases in (b) and (d). The bottleneck is more apparent in (b) than in (d), where the changes of  $f_{20,1}(t)/\theta$  stay within a five percent margin. This is due to slower rescaled time in the first scenario (a), where the population size drop affects polymorphisms prolongedly.



1 The third important feature of our model is the underlying assumption regarding genetic drift.  
2 The theoretical result for our SFS computations (Equation 33 in [15], see also SI 4) is based on  
3 a continuous time approximation of the Wright-Fisher model assuming descendants picking their  
4 parents at random from non-overlapping discrete generations. However, in our model, there is an  
5 overlap of generations in the host and parasite populations to allow the transmission of disease.  
6 A higher overlap of generations occurs for lower values of the death rate  $d$ . In order to apply the  
7 theoretical results, we therefore focus specifically on scenarios where the death rate  $d$  is close to the  
8 birth rate  $b = 1$ , so that most of the population is replaced within a 'generation' as defined by  $1/d$ .  
9 In addition, we build a simulation method that incorporates the overlap of generations in hosts and  
10 pathogens during the drawing of the next generation's allele frequencies.

## 11 **Modeling overlapping generations**

12 A simulation protocol was designed that accounts for the effect of overlapping generations in the  
13 stochastic sampling of the host and parasite SFS. The differential equations (1) and (2) are dis-  
14 cretized by choosing a sufficiently small value for  $\Delta t$  (its infinitesimal version being  $dt$ ) ensuring  
15 that the coevolutionary dynamics match the numerical evaluations from *Mathematica*. At every  
16 discrete generation  $\tau$ , the current population size  $N_\tau$  consists of  $N_o(\tau) = (1 - d) N(\tau - 1)$  in-  
17 dividuals that did not die and therefore overlap and  $N_b = N(\tau) - N_o(\tau)$  newborns (*i.e.* newly  
18 infected hosts in case of the parasite). As before, we assume that the SFS of both populations is in  
19 equilibrium at start of the infection, *i.e.* the population SFS at time zero is given by  $f_x(0) = \theta/x$   
20 with  $\theta = 2 N_{\text{ref}} \mu$  and  $\mu$  being the per generation mutation rate of an entire genome. The SFS of  
21 each population is recursively evaluated as follows. For a fixed generation  $\tau_0$ , the  $N(\tau_0)$  alleles  
22 are sampled from the pool of size  $N(\tau_0 - 1)$  and according to the allele frequencies at generation  
23  $\tau_0 - 1$ . The newborns and the overlap fraction are, respectively, obtained via sampling with and  
24 without replacement. New mutants arising only in newborns and as a single copy at a previously  
25 monomorphic site are obtained per generation as a Poisson random variable with mean  $2N_b \mu$  and  
26 added to the singleton class of the SFS. Once the population SFS  $f_x(t)$  is computed for a given time  
27 interval, its sample version  $f_{n,i}(t)$  is readily obtained via binomial sampling. Note that the number  
28 of novel mutations and therefore the amount of polymorphism over time is reduced by definition in  
29 this model with overlap compared to the Wright-Fisher model, where all individuals are newborns  
30 when descendants replace their parental generation.

## Results

1

### 2 **Effect of the various parameters on the dynamical system**

3 We are only interested in situations where at least one host and a parasite genotype survive and both  
4 populations coexist. Therefore, we derive first the criteria for the disease to spread in the population  
5 via the reproduction ratios (SI 2). Then, we scan the parameter space of our epidemiological  
6 model to determine the behavior of the coevolutionary dynamics and the speed of cycling. We thus  
7 eliminate situations where the disease spreads but the host and parasite populations finally collapse  
8 and get extinct or rise jointly in size. The behavior of the allele frequency cycling is determined  
9 by the state of the fixed point (computed in SI 3 for  $A = 2$ ). The cycling can be stable (regular  
10 cycles as fluctuating selection) or damp off to the fixed polymorphic state. The stability behavior  
11 of the hosts and the parasites, whose frequencies are not explicitly given by (1) and (2), cannot  
12 be explicitly determined by means of a Jacobian matrix and only be determined by numerically  
13 solving (1) and (2).

14 For the stability analysis, the initial conditions were chosen so that  $N^W$  is close to 10,000 and the  
15 infected alleles make up 20% of the healthy ones ( $H_1 = H_2 = 4150, I_{11} = I_{12} = I_{21} = I_{22} = 415$ ).  
16 The birth rates  $b$  were fixed to one and  $c_H = c_P = 0.05$  for the MA model and  $c_{H_1} = c_{P_2} =$   
17  $0.05, c_{H_2} = c_{P_1} = 0$  for the GFG model (see [20, 21] for comparable costs). The remaining  
18 parameters are given in the color-coded figures of SI 5 summarizing the results of the stability  
19 analysis, which become more diverse in behavior with a decreasing selection coefficient  $s$ , and an  
20 increasing difference between the birth rates  $b$  and the death rates  $d$ . We chose an appreciable  
21 number of initially infected alleles to omit cases where such genotypes could be instantaneously  
22 lost. Note that changing the birth rate to values other than one and employing various initial allele  
23 frequencies of healthy and infected alleles give results that correspond to the ones presented here  
24 (up to a rescaling of the underlying parameters).

25 The impact of the various parameters on the behavior of the dynamical system can be sum-  
26 marized as follows. An increasing difference between the birth and the death rates, *i.e.* between  
27  $b(1 - c_H)$  and  $d$ , results in 1) wider parameter ranges of the mortality  $\delta$  and the disease transmission  
28 rate  $\beta$  for which cycles may occur, and 2) also increases the number of cycles over a given time  
29 interval. While the number of cycles generally increases with increasing values of  $\delta$ ,  $\beta$  affects the  
30 number of cycles most distinctly for  $s = 1$ , for which smaller values of  $\beta$  lead to a reduced number  
31 of cycles. The speed of cycling for MA and GFG models is equivalent, if costs are set to zero. Thus,

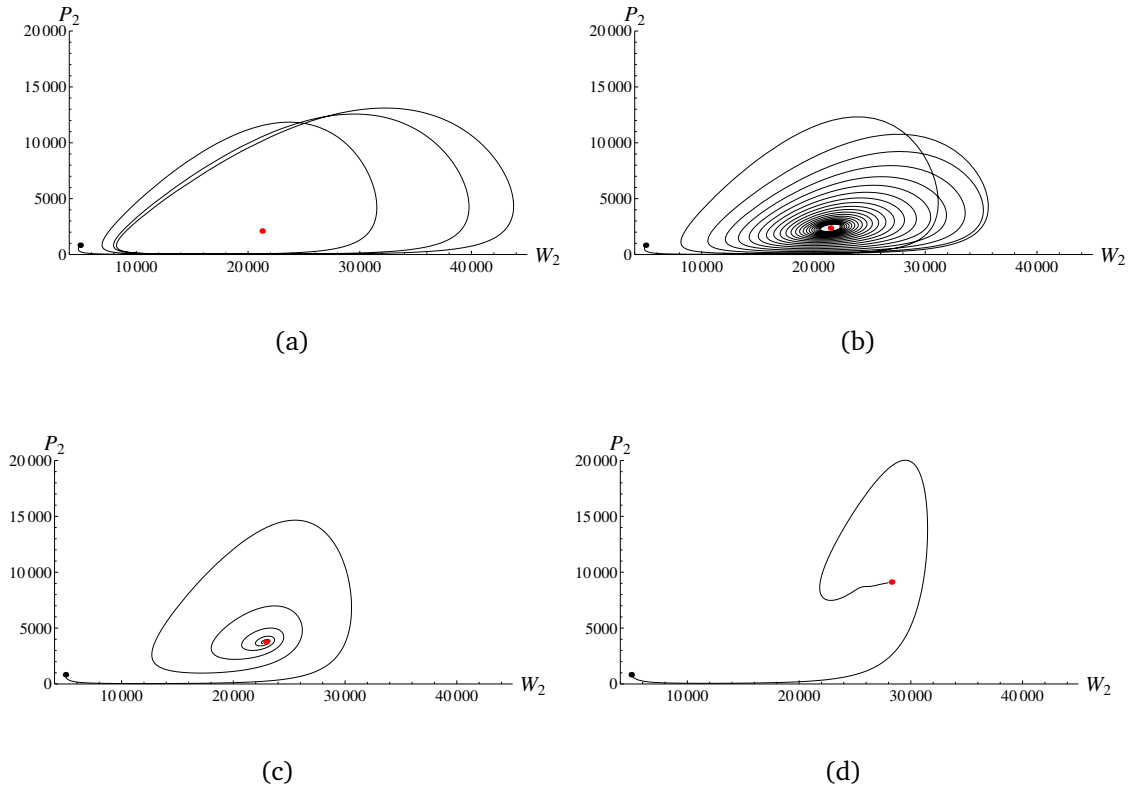


Figure 2 Parasite and host alleles of genotype two are plotted against each other for the GFG model over time by numerically solving (1) and (2) for the following parameter values (being equivalent for both genotypes):  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $\beta = 0.00005$ ,  $c_{H_1} = c_{P_2} = 0.05$  and  $c_{H_2} = c_{P_1} = 0$ . The initial conditions are  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$ . The selection coefficients  $s$  are given by (a) 1, (b) 0.9, (c) 0.6 and (d) 0.3. The parametric plots are shown for (a) 100, (b) 500, (c) 120 and (d) 96 time steps, which are the minimum amounts of time to complete one orbit of the final limit cycle (a), come close to the fixed point (b), (c), or to reach the fixed point (d). The initial and fixed points are colored in black and red, respectively.

1 it may be unrealistic to aim to infer the model itself (MA or GFG) based on polymorphism data. In  
 2 the following, we focus on the GFG model and present according results for the MA model in the  
 3 supplement.

4 As illustrated in Figure 2, besides the difference between  $b(1-c_H)$  and  $d$ , the selection coefficient  
 5  $s$  is the crucial parameter that determines the number of cycles per unit of time. A limit cycle is  
 6 observed for  $s = 1$ , a case defined as "castrating parasite". The cycling appears faster and with  
 7 quicker damping off with smaller values of  $s$  [2]. Note that employing the same parameters to the  
 8 MA model results in a loss of all parasite alleles for  $s = 1$  and  $s = 0.9$ , whereas cycling towards  
 9 the fixed point is obtained for  $s = 0.6$  and  $s = 0.3$ . The MA model with a death rate  $d = 0.6$   
 10 and otherwise equivalent parameters yields cycles for  $s = 1$  and  $s = 0.9$  as well (in contrast

1 to the behavior when  $d = 0.9$ ). The speed of cycling is therefore not only determined by the  
2 coevolutionary ( $s, c_H, c_P, \alpha$ ) and epidemiological parameters ( $\beta, \delta$ ), but also by the ecological  
3 characteristics of the species and the environment controlling the birth ( $b$ ) and death ( $d$ ) rates. An  
4 example of a MA model with a death rate  $d = 0.9$  and slow cycles that will be studied alongside the  
5 GFG example of Figure 2 is given in SI 6.

## 6 **Analyzing the polymorphism patterns of hosts and parasites**

### 7 *One parasite per host and equal generation times*

8 We provide detailed results for a few examples of GFG and MA models to highlight the key  
9 features regarding changes in the host and parasite SFS. Mutation rate and genetic drift occur on  
10 different time scales for hosts and parasites. The population size changes in the parasite occur  
11 on a relatively slower time scale compared to the host (Figure 3). Moreover, the amplitude of  
12 population size fluctuations is much more pronounced than in the host, the number of infected  
13 alleles fluctuating between about twenty and 13.000 individuals. The relatively weak and fast  
14 fluctuations of the host cannot be observed in the SFS except for a slight increase in the number of  
15 singletons over time. In contrast, the strong and relatively slow changes in the parasite population  
16 size are clearly reflected by the same statistic. The number of singletons decreases first due to the  
17 initial decline in the population size before tending to increase over time. A similar result for the  
18 MA model is given in SI 7.

19 We also observed that signatures of coevolutionary cycles in the host and parasite SFS de-  
20 pend mostly on the speed of their fluctuations in terms of population size scaled generation times,  
21 whereas the magnitude of the population size changes has a small influence on the SFS (rather  
22 apparent in Figure 1 than in Figure 3). We also evaluated one of the slowest cycling examples for  
23 a death rate of  $d = 0.3$  shown in the first panel in Figure SI 5.2.2 to illustrate that despite the  
24 difference in time scale in hosts and parasite cycles are also barely detectable in the parasite SFS  
25 (SI 8).

26 To detect changes in the SFS in host and parasite, it is also important to assess if enough genetic  
27 diversity can be observed over time. We find that the absolute number of polymorphisms strongly  
28 decreases over the considered time interval. For example in the GFG model (Figure 3) the number  
29 of segregating sites in the parasite decreases to about five percent and for the MA model (SI 7) even  
30 down to about one percent of their initial values. We also contrast two scenarios: 1) an initial total

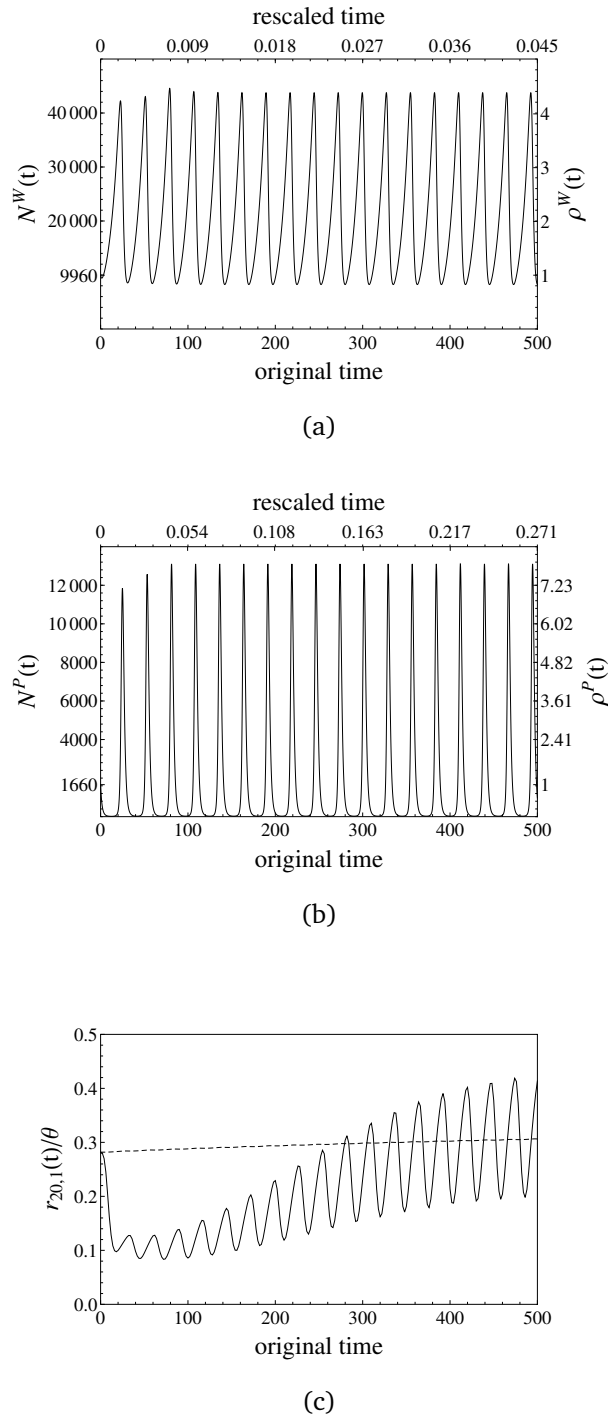


Figure 3 Population size changes in the host (a) and in the parasite (b) are generated for the GFG model via the parameters  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $\beta = 0.00005$ ,  $c_{H_1} = c_{P_2} = 0.05$ ,  $c_{H_2} = c_{P_1} = 0$ , and  $s = 1$ . The initial conditions are  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$ , so that the reference population sizes  $N_{\text{ref}}$  of the host and the parasite are given by 9960 and 1660, respectively, before their interaction starts at time zero. In both cases, the lower x-axes show time at the original scale of the dynamical system, whereas time is scaled by the respective values of  $N_{\text{ref}} \cdot 1/d$  for the upper x-axes. The left y-axes denote the absolute values of the changing population size and the right y-axes denote the relative values as  $\rho(t) = N(t)/N_{\text{ref}}$ . We evaluated the SFS for every second generation (on the original scale) in both cases and plot (c) the relative number of singletons  $r_{20,1}(t)$  against time for the host (dashed) and the parasite (solid).

Initial prevalence, $N_{\text{ref}}$ parasites	Fungi/Trypanosoma/Nematodes $\mu = 1$	Viruses $\mu = 0.1$	Bacteria $\mu = 0.001$
2%, 2000	[8232, 14191]	[823, 1419]	[8; 14]
20%, 1660	[560; 11779]	[56; 1178]	[1; 12]
50%, 3332	[48; 23642]	[5; 2364]	[0; 24]

Table 2 [minimum, maximum] of the number of segregating sites in full genome sequences of parasites following the GFG model and depending on the initial prevalence and population size. The minima are determined by the results in SI 9 and the maxima are the initial numbers of segregating sites. The per genome mutation rate,  $\mu$ , is an approximation based on genome length and per site mutation rate from different typical estimates [22].

1 host population size of 100,000 and 2% of initial disease prevalence, and 2) an initial population  
2 size of 10,000 with 50% initial prevalence (for GFG and MA models in SI 9). We observe that the  
3 initial prevalence defining the parasite population size plays a crucial role. Scenario two shows  
4 stronger fluctuations in the relative number of parasite singletons and a more drastic decrease in  
5 the absolute number of segregating sites over time than under scenario one. This shows that larger  
6 fluctuations in the relative SFS over time, which are in principle detectable in time sample polymor-  
7 phism data, go along with a stronger decrease in the total amount of observed polymorphisms (due  
8 to more drastic population bottlenecks). As a guideline, we provide an estimate of the possibility to  
9 detect changes in the SFS based on a sufficient number of segregating sites (the numerical minima  
10 and maxima in Figure SI 9.1 yield Table 2). Cycles are more likely observed in parasites with large  
11 genome mutation rates such as fungi or protozoans in contrast to bacteria (Table 2).

12 When comparing our computationally advantageous approach based on the Wright-Fisher  
13 model with our stochastic simulations, we find that polymorphism signatures as exemplarily  
14 measured by  $\Pi_{20}(t)$  agree in general for host and parasite over time (Figure 4). The parasite  
15 sample shows less polymorphism in the simulations with overlap while for the host this difference  
16 is negligible. The net effect of generation overlap under strong population bottlenecks with more  
17 dominant decline than expansion phases is an even stronger decrease of the effective population  
18 size and thus the amount of polymorphism, as seen in the parasite. This is due to less newborns  
19 contributing fewer novel mutations in the model with overlap and the different sampling schemes

1 of newborn and overlapping individuals. Whenever population size decreases, new offspring indi-  
2 viduals are present in smaller proportions than overlapping ones, whereas the reverse occurs when  
3 population size increases. The fraction of new offspring follows a sampling with replacement as for  
4 the Wright-Fisher model, while the overlap fraction is drawn without replacement. The difference  
5 of these two sampling schemes becomes apparent during phases of small population sizes. During  
6 decline (expansion) phases, the increased fraction of overlapping (newborn) individuals leads to  
7 stronger (lesser) deviations from the Wright-Fisher expectations. Consequently, as the parasite  
8 population is experiencing a drastic population decrease over time and several cycles, the amount  
9 of diversity differs between both approaches in contrast to the host population experiencing a  
10 slight increase in size over time.

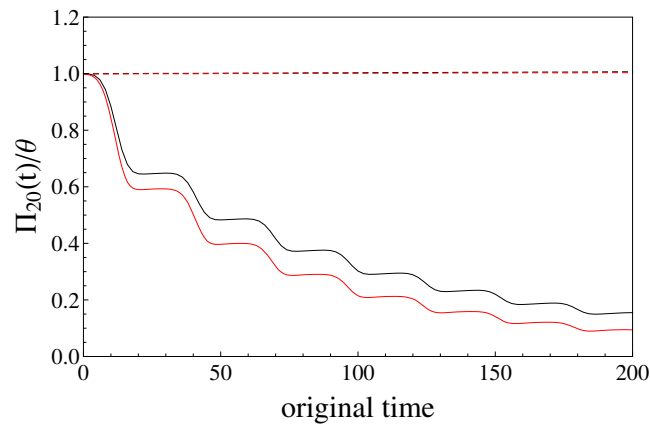


Figure 4 Population size changes in host (dashed) and parasite (solid) are generated for the GFG model with the same parameter set as in Figure 3. The average number of pairwise differences scaled by  $\theta$ ,  $\Pi_{20}(t)/\theta$ , is plotted against time for our analytical framework based on the diffusion approximation of the Wright-Fisher model (black) and for our simulation approach (red) with a fractional overlap of  $1 - d$  ( $d = 0.9$ ) between successive time steps. For the simulations, computational time steps are set to 0.001, a genome-wide mutation rate of  $\mu = 1$  is applied, and each value is obtained as an average over 10 repetitions.

### 11 *Multiple parasites per host and polycyclic diseases*

12 We extend our predictions for two classic deviations from our model. First, multiple or even  
13 a large number of parasites, as denoted by  $F$ , often infect a single host. Considering this effect  
14 on the SFS leads to an increased scaled mutation rate  $F\theta$  thereby increasing the number of seg-  
15 regating sites that can be detected by full-genome sequencing. Concurrently, relative time is sped  
16 up by  $F$  due to the larger initial (reference) population size of the parasite. Therefore, two op-  
17 posite effects are expected when increasing the number of parasites per host, an increase in the

1 amount of polymorphism comes at the cost of a reduced amount of detectable cycles. Second, host  
2 and parasite generation times may differ from one another with parasites often exhibiting smaller  
3 generation times especially for virus, bacteria or fungi (polycyclic diseases). We define  $E$  parasite  
4 generations per host generation so that the relative time for the parasite is slowed down by  $E$ . This  
5 rescaling is expected to enhance the detectability of coevolutionary cycles using the parasite SFS.  
6 We investigate the joint impact of multiple parasites per host and polycyclic diseases in SI 10.

7 We compute the SFS over time for the GFG above (Figure 3) with nine different combinations  
8 of values for  $E$  and  $F$  (see SI 10). To evaluate fluctuations in the polymorphism pattern over  
9 time, we compute the relative rate of change of  $\Pi_{20}(t)$  (SI 10) at various equidistantly distributed  
10 sampling points over time. These points are chosen to be fixed in time and independent of parasite  
11 generation time, so that the various cases are equivalently clocked on the original time scale of the  
12 dynamical system. As illustrated in SI 10 more sampling points are needed to capture the cycling  
13 for multiple parasites per host ( $F$ ), whereas less samples are needed for polycyclic diseases. We  
14 therefore show here that the number of time samples necessary to recover the number of cycles  
15 can be determined knowing the biology of host and parasite.

## 16 Discussion

17 We develop here a model to analyze the evolution of neutral genome-wide polymorphism of coe-  
18 volving host and parasite populations. Variations in polymorphism reflect the co-demographic his-  
19 tory driven by eco-evolutionary feedbacks between 1) the frequency changes in host resistance and  
20 parasite infectivity over time due to frequency-dependent selection, and 2) the ecological changes  
21 in host and parasite population sizes. While antagonistic or synergistic coevolution is a process  
22 driven by natural selection, our approach is the first to provide a description of the consequences  
23 of coevolutionary dynamics on neutral genome-wide polymorphism.

24 We demonstrate that using time sampling data, *i.e.* population samples of hosts and parasites at  
25 different time points, it is possible to track the existence and speed of eco-evolutionary cycles using  
26 polymorphism data. Our main result states that cycles of coevolution are detectable in the parasite  
27 and barely in the host population. This is due to a fundamental difference in the time scale of  
28 neutral evolution between interacting species and the strength of their size fluctuations over time  
29 both crucially depending on the initial population size at the onset of epidemics (*i.e.* the start of  
30 the coevolutionary history). The time scale of neutral evolution is also determined by the parasite



1 generation time and number of parasites per infected host. We study here one coevolutionary run  
2 starting by the introduction of a parasite population into a larger host population and generating  
3 a dynamics over several hundreds of generations. If new infectivity or resistance alleles appear by  
4 mutation, the epidemics and the cycling behavior are affected, and our model should be reset to  
5 evaluate a new run. The time scale that we investigate is therefore intermediate between the classic  
6 expectations of long coevolution and its signatures at interacting loci [6, 23] and the short-term  
7 epidemiological dynamics (with susceptible hosts and one parasite type) [24].

8 Polymorphism data can be used to detect coevolutionary cycles, if such cycles run sufficiently  
9 long at adequately low speed. We indeed predict that long term occurrence of cycles should be  
10 searched for in parasites that strongly decrease the host fitness due to high disease severity  $s$  (par-  
11 asite effect on fecundity). For low to moderate disease transmission and smaller values of  $s$  the  
12 internal polymorphic equilibrium point is a stable attractor, meaning that cycles damp off quickly  
13 towards a polymorphic equilibrium at which population size and allele frequencies are fixed (Fig-  
14 ure 2). For high values of the disease transmission rate  $\beta$  and parasite virulence  $\delta$  (effect of parasite  
15 on mortality), the internal polymorphic equilibrium point is an unstable point [2], and a monomor-  
16 phic equilibrium is reached with a fixed population size [10, 12]. Cycles should be slow enough to  
17 be observable in polymorphism data, and our results challenge the classic assumption that coevolu-  
18 tionary cycles are too fast for being observed in the SFS. Interestingly, the speed of cycling depends  
19 mainly on two ecological parameters, *i.e.* the birth rate  $b$  and the death rate  $d$  of the host, and to a  
20 lesser extent on the coevolutionary parameters ( $s$ , costs of resistance and virulence).

21 Eco-evolutionary cycles occur and are observable in polymorphism data, when the effect of the  
22 parasite on host fecundity ( $s$ ) is sufficiently strong. We predict that our results are applicable to  
23 many host-parasite systems with castrating parasites, whose transmission is associated with host  
24 death (algae-rotifer [25], bacteria-phage [26] and *Daphnia magna*-bacterium [27]). For plant  
25 pathogens, cycles may be less observable because the disease severity can range from low to very  
26 high (or even castrating, [28]), but often depends on abiotic factors [29]. Most epidemiological  
27 studies have focused on the evolution of virulence (effect of parasite on host mortality) and disease  
28 transmission within the short duration of an epidemics. However, to use polymorphism data for the  
29 study of eco-evolutionary dynamics, parasite virulence is not an essential parameter to be measured  
30 or estimated, as it is more useful to quantify the difference between the host's birth and natural  
31 death rates. Another practical reason for favoring hosts with high death rates is that our SFS  
32 computations are based on the Wright-Fisher model assuming non-overlapping generations [15].

1 Note that in epidemiological models [2, 10, 12] overlapping generations in the host are a necessary  
2 assumption, since a disease can only be transmitted among living hosts. We wrote a code in C  
3 (available from the authors upon request) that explicitly accounts for less newborns contributing  
4 fewer mutations while generations overlap (and compared to the Wright-Fisher model) to evaluate  
5 the SFS. The simulations show that our Wright-Fisher approximation is robust with respect to  
6 overlapping generations for hosts with high death rates.

7 We only consider neutral sites because arbitrary demographic changes can be used in the an-  
8 alytic solution for the SFS as needed to cope with the complex demographics arising from our  
9 dynamical system. The frequency spectrum for sites under selection can only be computed for  
10 piecewise changes in the population size [30]. This approach is not readily applicable for such  
11 complex demographics because their discretization would be computationally cumbersome.

12 Time sampling is crucial for capturing cycles (see SI 10) that can be observed in the poly-  
13 morphism data for several hundreds of parasite generations, if the genome mutation rate and the  
14 effective population size are sufficiently large (see Table 2) and if SNPs are sampled at appropriate  
15 time points (see SI 10). To test our predictions, time samples can be readily obtained in experi-  
16 mental coevolution set-ups [25, 26], whereas this may be more complex for natural populations.  
17 Nevertheless, samples from the past can be obtained for crustaceans (*Daphnia*, [27]) from dormant  
18 stages deposited in sediments, and for plant species from seeds in the soil (possibly using ancient  
19 DNA recovery techniques).

20 We predict that parasites undergoing several generations per host generation but producing  
21 small amount of pathogen propagules per host should be the species exhibiting most clearly the  
22 signature of co-demographic dynamics in polymorphism data. Our results pave the way to use time  
23 sample genomic data of hosts and parasites from wild or experimental populations, to analyze,  
24 infer and take into account the co-demographic history of the antagonistic species and scan for  
25 genes under coevolution with greater accuracy.

## 26 Funding

27 DZ and SJ were supported by the Deutsche Forschungsgemeinschaft (DFG) grant STE325/14 from  
28 the Priority Program SPP1590 "Probabilistic Structures in Evolution" to WS. MV was supported  
29 by DFG grant TE809/1 to AT, and SJ by grant TE809/3 from the SPP1819 "Rapid Evolutionary  
30 Adaptation" to AT.

## References

- 1 [1] S. Gandon, T. Day, C. J. E. Metcalf, B. T. Grenfell, Forecasting epidemiological and evolution-  
2 ary dynamics of infectious diseases, *Trends in Ecology & Evolution* 31 (2016) 776–788.
- 3 [2] B. Ashby, M. Boots, Multi-mode fluctuating selection in host–parasite coevolution, *Ecology*  
4 *Letters* 20 (2017) 357–365.
- 5 [3] B. Clarke, P. O’Donald, Frequency-dependent selection, *Heredity* 19 (1964) 201–206.
- 6 [4] A. Tellier, J. K. Brown, Stability of genetic polymorphism in host–parasite interactions, *Pro-*  
7 *ceedings of the Royal Society of London B: Biological Sciences* 274 (2007) 809–817.
- 8 [5] J. Bergelson, M. Kreitman, E. A. Stahl, D. Tian, Evolutionary dynamics of plant R-genes,  
9 *Science* 292 (2001) 2281–2285.
- 10 [6] M. E. Woolhouse, J. P. Webster, E. Domingo, B. Charlesworth, B. R. Levin, Biological and  
11 biomedical implications of the co-evolution of pathogens and their hosts, *Nature Genetics* 32  
12 (2002) 569–577.
- 13 [7] E. A. Stahl, G. Dwyer, R. Mauricio, M. Kreitman, J. Bergelson, Dynamics of disease resistance  
14 polymorphism at the *Rpm1* locus of *Arabidopsis*, *Nature* 400 (1999) 667–671.
- 15 [8] J. K. Brown, A. Tellier, Plant-parasite coevolution: bridging the gap between genetics and  
16 ecology, *Annual Review of Phytopathology* 49 (2011) 345–367.
- 17 [9] R. M. May, R. Anderson, Epidemiology and genetics in the coevolution of parasites and hosts,  
18 *Proceedings of the Royal Society of London B: Biological Sciences* 219 (1983) 281–313.
- 19 [10] A. Tellier, J. K. M. Brown, The Influence of Perenniality and Seed Banks on Polymorphism in  
20 Plant-Parasite Interactions, *American Naturalist* 174 (2009) 769–779.
- 21 [11] S. A. Frank, Coevolutionary genetics of hosts and parasites with quantitative inheritance,  
22 *Evolutionary Ecology* 8 (1994) 74–94.
- 23 [12] C. S. Gokhale, A. Papkou, A. Traulsen, H. Schulenburg, Lotka–Volterra dynamics kills the Red  
24 Queen: population size fluctuations and associated stochasticity dramatically change host-  
25 parasite coevolution, *BMC Evolutionary Biology* 13 (2013) 254.
- 26

- 1 [13] Y. Song, C. S. Gokhale, A. Papkou, H. Schulenburg, A. Traulsen, Host-parasite coevolution in  
2 populations of constant and variable size, *BMC evolutionary biology* 15 (2015) 212.
- 3 [14] J. F. Rabajante, J. M. Tubay, H. Ito, T. Uehara, S. Kakishima, S. Morita, J. Yoshimura, D. Ebert,  
4 Host-parasite Red Queen dynamics with phase-locked rare genotypes, *Science Advances* 2  
5 (2016) e1501548.
- 6 [15] D. Živković, W. Stephan, Analytical results on the neutral non-equilibrium allele frequency  
7 spectrum based on diffusion theory, *Theoretical Population Biology* 79 (2011) 184–191.
- 8 [16] M. Boots, A. White, A. Best, R. Bowers, How specificity and epidemiology drive the coevolu-  
9 tion of static trait diversity in hosts and parasites, *Evolution* 68 (2014) 1594–1606.
- 10 [17] R. C. Griffiths, S. Tavaré, The age of a mutation in a general coalescent tree., *Stochastic*  
11 *Models* 14 (1998) 273–295.
- 12 [18] D. Živković, T. Wiehe, Second-order moments of segregating sites under variable population  
13 size., *Genetics* 180 (2008) 341–357.
- 14 [19] S. Gandon, Y. Michalakis, Local adaptation, evolutionary potential and host–parasite coevolu-  
15 tion: interactions between migration, mutation, population size and generation time, *Journal*  
16 *of Evolutionary Biology* 15 (2002) 451–462.
- 17 [20] P. H. Thrall, J. J. Burdon, Evolution of virulence in a plant host-pathogen metapopulation,  
18 *Science* 299 (2003) 1735–1737.
- 19 [21] D. Tian, M. B. Traw, J. Q. Chen, M. Kreitman, J. Bergelson, Fitness costs of R-gene-mediated  
20 resistance in *Arabidopsis thaliana*, *Nature* 423 (2003) 74–77.
- 21 [22] M. Lynch, M. S. Ackerman, J.-F. Gout, H. Long, W. Sung, W. K. Thomas, P. L. Foster, Genetic  
22 drift, selection and the evolution of the mutation rate., *Nature Reviews Genetics* 17 (2016)  
23 704–714.
- 24 [23] A. Tellier, S. Moreno-Gómez, W. Stephan, Speed of adaptation and genomic footprints of  
25 host–parasite coevolution under arms race and trench warfare dynamics, *Evolution* 68 (2014)  
26 2211–2224.
- 27 [24] T. Stadler, R. Kouyos, V. von Wyl, S. Yerly, Böni, J., *et al.*, Estimating the basic reproductive  
28 number from viral sequence data, *Molecular Biology and Evolution* 29 (2012) 347–357.

- 1 [25] L. Becks, S. P. Ellner, L. E. Jones, N. G. Hairston, The functional genomics of an eco-  
2 evolutionary feedback loop: linking gene expression, trait evolution, and community dynam-  
3 ics, *Ecology Letters* 15 (2012) 492–501.
- 4 [26] A. R. Hall, P. D. Scanlan, A. D. Morgan, A. Buckling, Host–parasite coevolutionary arms races  
5 give way to fluctuating selection, *Ecology Letters* 14 (2011) 635–642.
- 6 [27] E. Decaestecker, S. Gaba, J. A. M. Raeymaekers, R. Stoks, L. Van Kerckhoven, D. Ebert, L. De  
7 Meester, Host–parasite ‘Red Queen’ dynamics archived in pond sediment., *Nature* 450 (2007)  
8 870–873.
- 9 [28] K. Clay, Parasitic castration of plants by fungi, *Trends in Ecology & Evolution* 6 (1991) 162–  
10 166.
- 11 [29] G. N. Agrios, *Plant Pathology*, 5th Edition, Academic Press, 2005.
- 12 [30] D. Živković, M. Steinrücken, Y. S. Song, W. Stephan, Transition densities and sample fre-  
13 quency spectra of diffusion processes with selection and variable population size., *Genetics*  
14 200 (2015) 601–617.

## Supplementary Information

### Neutral genomic signatures of host-parasite coevolution

Daniel Živković, Sona John, Mélissa Verin, Wolfgang Stephan, Aurélien Tellier

.

## 1 SI 1. The host effective population size over time

2 From (1) and (2) we immediately obtain

$$\frac{dN^W}{dt} = \sum_{i=1}^A H_i [b_i(1 - c_{H_i}) - d_i] + \sum_{i=1}^A b_i(1 - c_{H_i}) \sum_{j=1}^A (1 - s_{ij}) I_{ij} - \sum_{i=1}^A \sum_{j=1}^A (d_i + \delta_{ij}) I_{ij}.$$

3 Assuming that  $s_{ij} = s_i$ ,  $\delta_{ij} = \delta_i$  (i.e. both parameters being independent of the parasite genotype)

4 and setting  $dN^W/dt$  to zero, we have

$$\sum_{i=1}^A H_i [b_i(1 - c_{H_i}) - d_i] + \sum_{i=1}^A [b_i(1 - c_{H_i})(1 - s_i) - (d_i + \delta_i)] \sum_{j=1}^A I_{ij} = 0.$$

5 It particularly follows that  $s_i = \delta_i = 0$  requires  $d_i = b_i(1 - c_{H_i})$  to have a constant population size

6 in the host for arbitrary choices of  $H_i$  and  $I_{ij}$ .

## 7 SI 2. Basic reproduction ratios

8 From (2), we obtain

$$\frac{dP_j}{dt} = - \sum_{i=1}^A I_{ij}(d_i + \delta_{ij}) + H_i \left[ \alpha_{ij} \beta_{ij} (1 - c_{P_j}) \sum_{k=1}^A I_{kj} \right].$$

9 Assuming that  $d_i = d$  and  $\delta_{ij} = \delta_j$  (i.e. both parameters being independent of the host genotype),

10 the former equation simplifies to

$$\frac{dP_j}{dt} = P_j \left[ -(d + \delta_j) + \sum_{i=1}^A H_i \alpha_{ij} \beta_{ij} (1 - c_{P_j}) \right].$$

11 The reproduction ratios of the parasite genotypes are given by

$$R_{0,j} = \sum_{i=1}^A H_i \alpha_{ij} \beta_{ij} (1 - c_{P_j}) / (d + \delta_j).$$

12 If all  $R_{0,j} < 1$ , the parasite genotypes are eliminated since they kill more hosts than they infect new

13 healthy ones. When one of these ratios is greater than one, the corresponding parasite genotype is

14 maintained in the population.

### 1 **SI 3. Fixed points of the dynamical system**

2 We calculated the fixed points by setting (1) and (2) to zero. Note that the solutions were first  
 3 obtained in terms of the numbers of healthy and infected hosts before being added up to obtain the  
 4 fixed points of hosts and parasites assuming one parasite per host. The results for the two-allele  
 5 case are summarized below. For the MA and iGFG model the results for more than two alleles  
 6 correspond to those of two alleles, whereas results for the iMA and the GFG model can only be  
 7 obtained for special cases when  $A = 3$ , but the solutions are not enlightening.

$$\text{Define } u_i = \frac{(b_i(1 - c_{H_i})s_{ii} + \delta_{ii})(\delta_{ii} + d_i)}{\beta_{ii}(-b_i(1 - c_{H_i})(1 - s_{ii}) + \delta_{ii} + d_i)} \quad \text{and} \quad v_i = \frac{(b_i(1 - c_{H_i}) - d_i)(\delta_{ii} + d_i)}{\beta(-b_i(1 - c_{H_i})(1 - s_{ii}) + \delta_{ii} + d_i)}.$$

8 For the **MA model** the fixed point, where all alleles may have nonzero frequencies, is given by  
 9  $(W_1^*, W_2^*, P_1^*, P_2^*)$ , where for  $i = 1, 2$ ,

$$W_i^* = \frac{u_i}{1 - c_{P_i}} \quad \text{and} \quad P_i^* = \frac{v_i}{1 - c_{P_i}}.$$

10 Besides the trivial solution of all alleles having frequency zero, the remaining solutions are given  
 11 by  $(W_1^*, 0, P_1^*, 0)$  and  $(0, W_2^*, 0, P_2^*)$ .

12  
 13 For the **iMA model** the fixed point, where all alleles may have nonzero frequencies, is given by  
 14  $(W_1^{**}, W_2^{**}, P_1^{**}, P_2^{**})$ , where for  $i = 1, 2$ ,

$$W_i^{**} = \frac{u_i}{1 - c_{P_{3-i}}} \quad \text{and} \quad P_{3-i}^{**} = \frac{v_i}{1 - c_{P_{3-i}}}.$$

15 Besides all alleles having frequency zero, the remaining solutions are given by  $(W_1^{**}, 0, 0, P_2^{**})$  and  
 16  $(0, W_2^{**}, P_1^{**}, 0)$ .

17  
 18 For the **GFG model** the fixed point, where all alleles may have nonzero frequencies, can be  
 19 evaluated but is of complicated form. So we only note the results for  $\beta_{ij} = \beta_i$ ,  $\delta_{ij} = \delta_i$  and  
 20  $s_{ij} = s_i$  (*i.e.* all of these parameters are independent of the parasite genotype). The fixed point



1 with nonzero entries is given by  $(W_1^{***}, W_2^{***}, P_1^{***}, P_2^{***})$ , where

$$W_1^{***} = \frac{(c_{P_2} - c_{P_1})u_1}{(1 - c_{P_2})(1 - c_{P_1})}, \quad W_2^{***} = \frac{u_2}{1 - c_{P_1}}, \quad P_1^{***} = \frac{v_2 - v_1}{1 - c_{P_1}} \quad \text{and} \quad P_2^{***} = \frac{v_1}{1 - c_{P_2}}.$$

2 Besides all alleles having frequency zero, further solutions are given by  $(0, W_2^*, 0, P_2^*)$ ,  $(W_1^{**}, 0, 0, P_2^{**})$

3 and  $(0, W_2^{**}, P_1^{**}, 0)$ .

4

5 For  $c_{H_1} = c_{H_2} = c_{P_1} = c_{P_2} = 0$ , we also obtain

$$W_1^{***} = 0, \quad W_2^{***} = -\hat{I}_{22} + u_2, \quad P_1^{***} = -\hat{I}_{22} + v_2 \quad \text{and} \quad P_2^{***} = 0,$$

6 where  $\hat{I}_{22}$  denotes the equilibrium solution of the infected host genotype. With the additional as-

7 sumption of equivalent rates and costs among both host and parasite genotypes, so that particularly

8  $u_1 = u_2 = u$  and  $v_1 = v_2 = v$ , we further have

$$W_1^{***} = u \left( 1 - \frac{\beta}{\delta + d} \hat{H}_2 \right), \quad W_2^{***} = u \frac{\beta}{\delta + d} \hat{H}_2, \quad P_1^{***} = v \left( 1 - \frac{\beta}{\delta + d} \hat{H}_2 \right) \quad \text{and} \quad P_2^{***} = v \frac{\beta}{\delta + d} \hat{H}_2,$$

9 where  $\hat{H}_2$  denotes the equilibrium solution of the second healthy host genotype.

10

11 For the **iGFG** model, the nontrivial equilibrium solution is equivalent to the MA model with

12  $(W_1^*, 0, P_1^*, 0)$ .

## 1 SI 4. Summary of analytical results for site frequency spectra and related statistics

2 Assume a model, where a haploid population is evolving according to Wright-Fisher dynamics  
 3 forwards in time and being of size  $N_{\text{ref}}$  at (and before) time zero. A sample of  $n$  sequences is  
 4 taken. Neutral mutations occur at unlinked and previously monomorphic sites at rate  $\theta = 2 N_{\text{ref}} \mu$ ,  
 5  $\mu$  being the mutation rate per genome per generation. Scale time in units of  $N_{\text{ref}}$  generations and  
 6 let  $N_{\text{ref}} \rightarrow \infty$  to reach the diffusion limit. Thereby, the relative population size  $N(t)/N_{\text{ref}}$  converges  
 7 to the strictly positive and piecewise continuous scaling function  $\rho(t)$ . The site frequency spectrum  
 8 is the distribution of the number of times a mutant allele is observed in the sample among the  
 9 polymorphic loci.

10

11 The absolute site-frequencies over time  $f_{n,i}(t)$ ,  $1 \leq i \leq n - 1$ , with a mutation-drift equilibrium at  
 12 time zero can be simply obtained from Equation 33 in [15] as

$$f_{n,i}(t) = \frac{\theta}{i} \sum_{k=2}^n (-1)^k (2k-1) {}_3F_2(n-i+1, k, 1-k; n+1, 2; 1) \left[ R(0) + \binom{k}{2} \int_0^t R(s) ds \right],$$

13 where  ${}_3F_2(a, b, c; d, e; z) = \sum_{l \geq 0} (a_{(l)} b_{(l)} c_{(l)}) / (d_{(l)} e_{(l)}) z^l / l!$ , with  $p_{(0)} = 1$  and  $p_{(l)} = p(p+1) \cdots (p+l-1)$

14 for  $l \geq 1$ , is a *generalized hypergeometric function*, and  $R(s) = \exp \left[ -\binom{k}{2} \int_s^t \rho^{-1}(u) du \right]$ .

15

16 The relative site frequencies over time  $r_{n,i}(t)$  are obtained as  $r_{n,i}(t) = f_{n,i}(t) / \sum_{k=1}^{n-1} f_{n,k}(t)$ , where  
 17 the denominator gives the absolute number of segregating sites,  $S_n(t)$ . For the average number of  
 18 pairwise differences, we have  $\Pi_n(t) = 1 / \binom{n}{2} \sum_{k=1}^{n-1} k(n-k) f_{n,k}(t)$ .

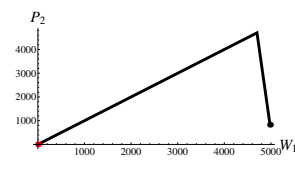
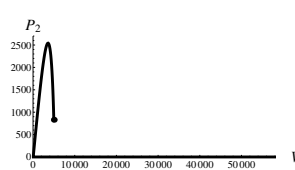
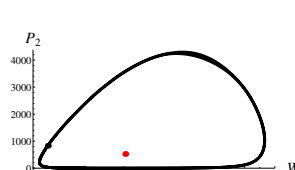

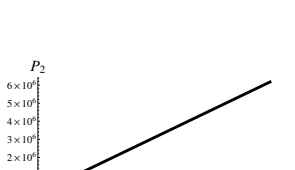
19

20 These implementations can be computationally demanding since the inverse of the relative popula-  
 21 tion size function  $\rho(t)$  has to be integrated numerically with high precision before being applied to  
 22 an exponential function that has to be integrated numerically again from the initial time zero up to  
 23 the time points at which the SFS are evaluated. These exponentials also cover binomial coefficients  
 24 that depend on sample size and the numerical integration has to be performed for every single one  
 25 of those. Therefore, an increasing amount of integration steps and computation time is required  
 26 with increasing sample size, so that we only consider a rather small sample size of twenty individ-  
 27 uals for which the SFS of a single time point can be rather quickly obtained and even for later time  
 28 points. This is particularly desirable, since we are interested in taking samples recurrently over the  
 29 course of time.

1 **SI 5. Stability behavior of the fixed points for two alleles**

2 *SI 5.1. The matching-allele model*

Table SI 5.1.1 Legend of coloring for possible allele frequency changes

color: type of allele frequency change	exemplary figure
black: loss of all alleles	
blue: $H_1$ and $H_2$ increase towards infinity; loss of all infected alleles	
orange: limit cycle; $I_{12}$ and $I_{21}$ get lost	
red: circular to rectilinear motion towards a stable equilibrium; $I_{12}$ and $I_{21}$ get lost;	
green: $H_1$ and $H_2$ converge to certain numbers; $I_{12}$ and $I_{21}$ get lost; $I_{11}$ and $I_{22}$ increase towards infinity	

In the left column the possible fates of (healthy and infected) genotypes are summarized and assigned to colors. In the right column corresponding exemplary parametric plots of host and parasite allele frequencies are shown using Equations (1) and (2). The black dot depicts the initial allele frequency, whereas the red dot shows the zero point in the first and the fixed point in the third and in the fourth subfigure.

Figure SI 5.1.1

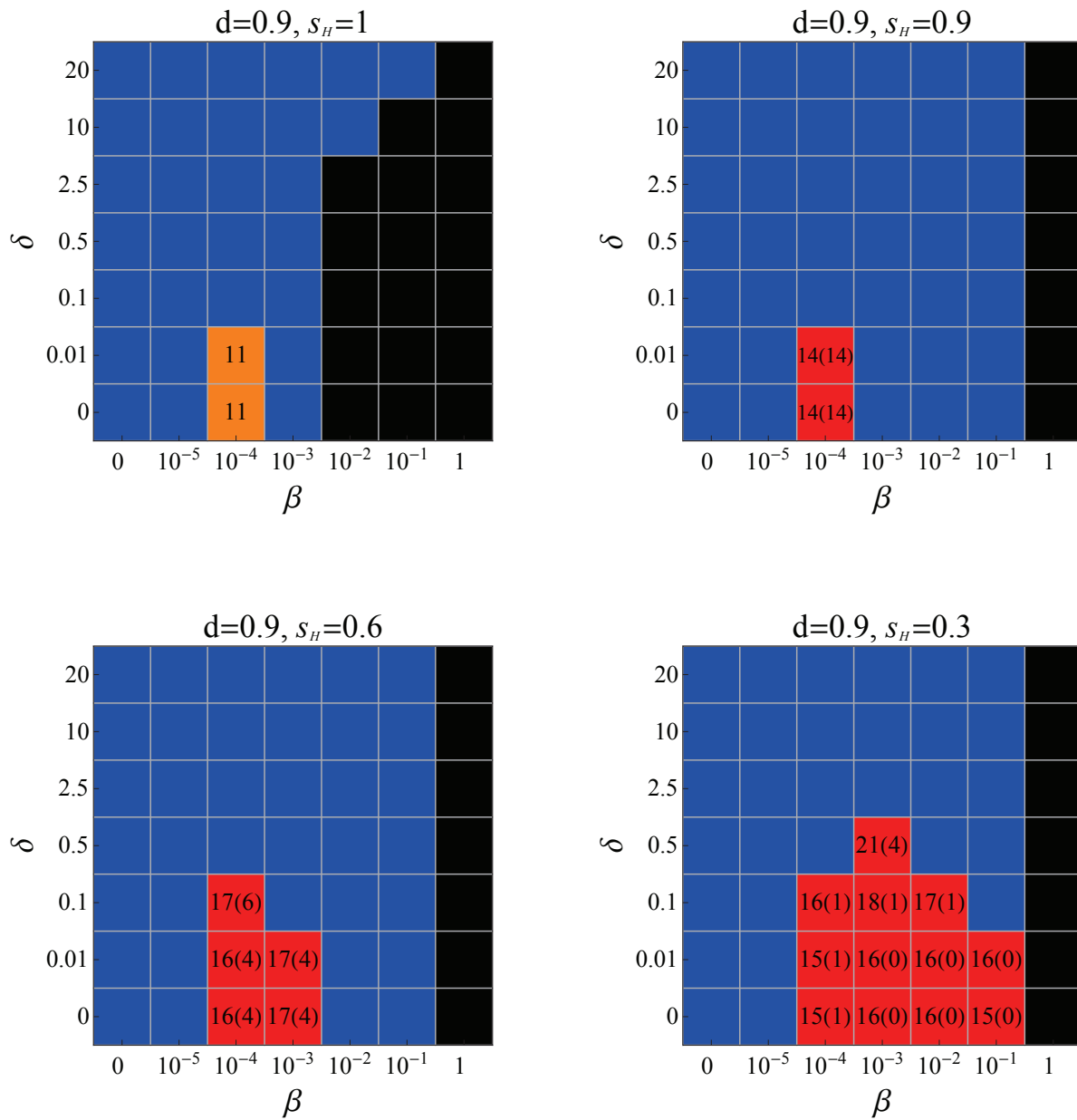
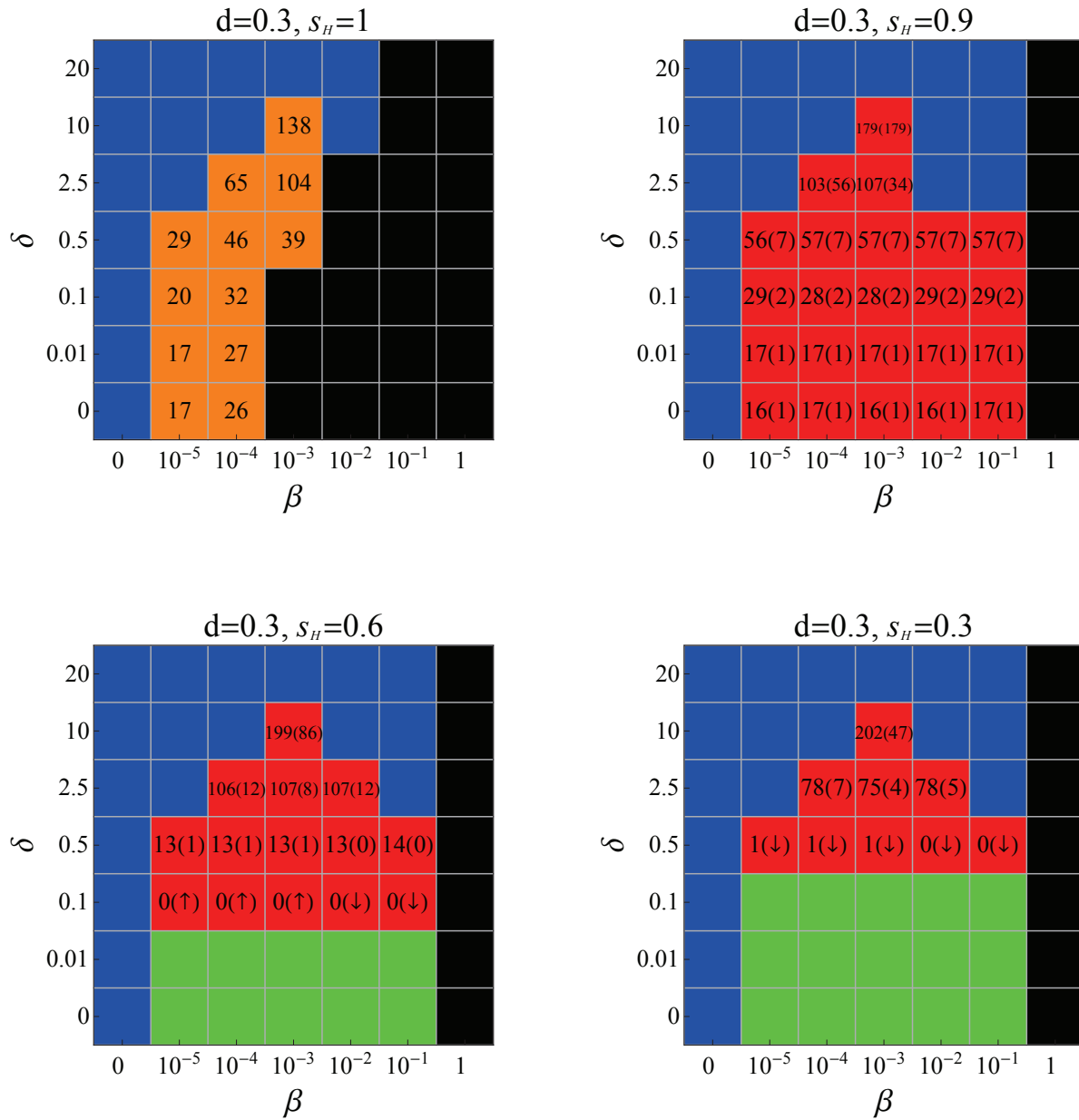


Figure SI 5.1.2

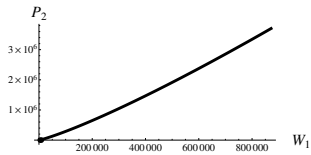
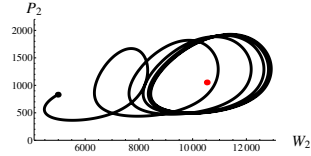
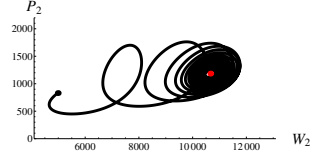
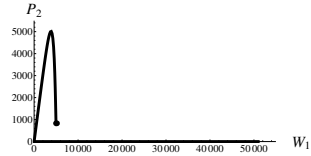
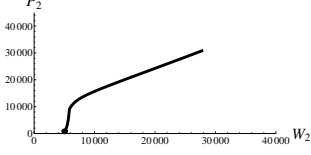


1 *Description of the results*

2 For each panel of Figures SI 5.1.1 and SI 5.1.2 the system of differential equations (1) and (2)  
3 were numerically evaluated for  $b = 1$ ,  $c_P = c_H = 0.05$  and the specified parameters over 500 time  
4 steps. The various possible coevolutionary scenarios of hosts and parasites are color-coded and  
5 summarized in detail in Table SI 5.1.1. The squares featuring numbers represent the parameter  
6 combinations for which cycling around or immediate attainment of the fixed point of  $N^W$  is ob-  
7 tained. For  $s = 1$ , the amplitude of the cycles around the fixed point of  $N^W$  remains constant over  
8 time. For other values of  $s$  we distinguish two cases: We first evaluated the total amount of cycles  
9 around the fixed point of  $N^W$  based on a numerical accuracy of 20 digits. Since the cycling dumps  
10 off in these cases, we also note in brackets the number of cycles until the amplitude is enclosed  
11 by the interval  $[0.95 N^W, 1.05 N^W]$ . The upwards and downwards arrows in some of the brackets  
12 illustrate expansions and declines to the fixed point, respectively.

1 SI 5.2. The gene-for-gene model

Table SI 5.2.1 Legend of coloring for possible allele frequency changes

color: type of allele frequency change	exemplary figure
<p>green: <math>H_1</math> and <math>H_2</math> converge to certain numbers  (<math>H_1</math> can get lost); <math>I_{11}</math> gets lost;  <math>I_{21}</math> and <math>I_{22}</math> increase towards infinity.  <math>I_{12}</math> increases towards infinity (dark green),  or gets lost (light green)</p>	
<p>orange: limit cycle; <math>H_1</math>, <math>I_{11}</math> and <math>I_{12}</math> get lost.  <math>I_{21}</math> (dark orange) or  <math>I_{22}</math> (light orange) gets lost</p>	
<p>red: circular to rectilinear motion  towards a stable equilibrium;  <math>H_1</math>, <math>I_{11}</math> and <math>I_{12}</math> get lost.  <math>I_{21}</math> (dark red) or  <math>I_{22}</math> (light red) gets lost</p>	
<p>yellow: <math>H_1</math> increases towards infinity;  loss of all the other alleles</p>	
<p>gray: <math>H_2</math> converges to a certain number;  <math>H_1</math>, <math>I_{11}</math> and <math>I_{12}</math> get lost.  <math>I_{22}</math> increases towards infinity and  <math>I_{21}</math> gets lost (dark gray),  or vice versa (light gray)</p>	

In the left column the possible fates of (healthy and infected) genotypes are summarized in addition or slightly modified to the scenarios already presented in Table SI 5.1.1 by means of the MA model. In the right column exemplary parametric plots of host and parasite allele frequencies are shown (in the first three panels for the dark color). The black dot depicts again the initial allele frequency.

Figure SI 5.2.1

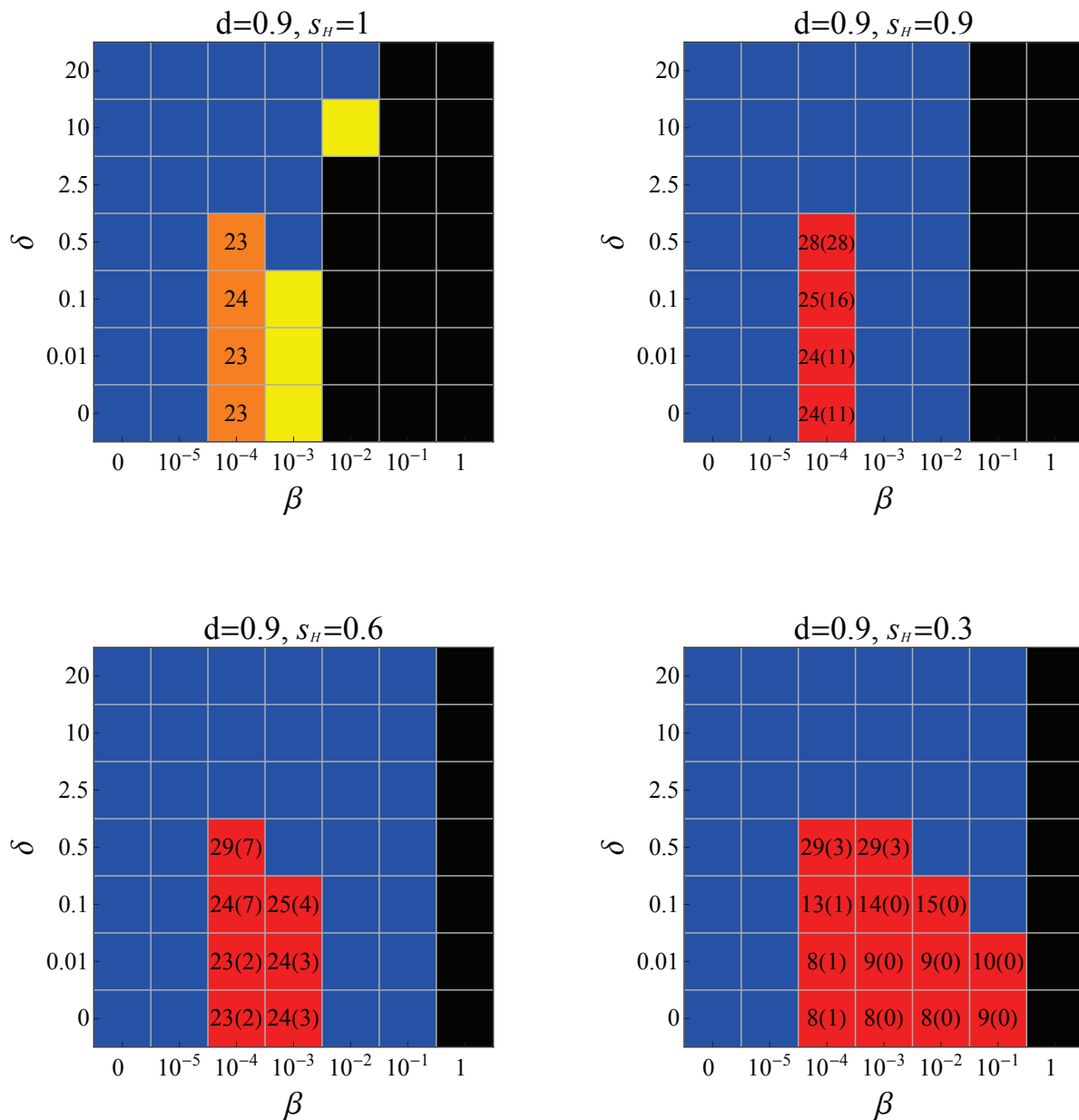
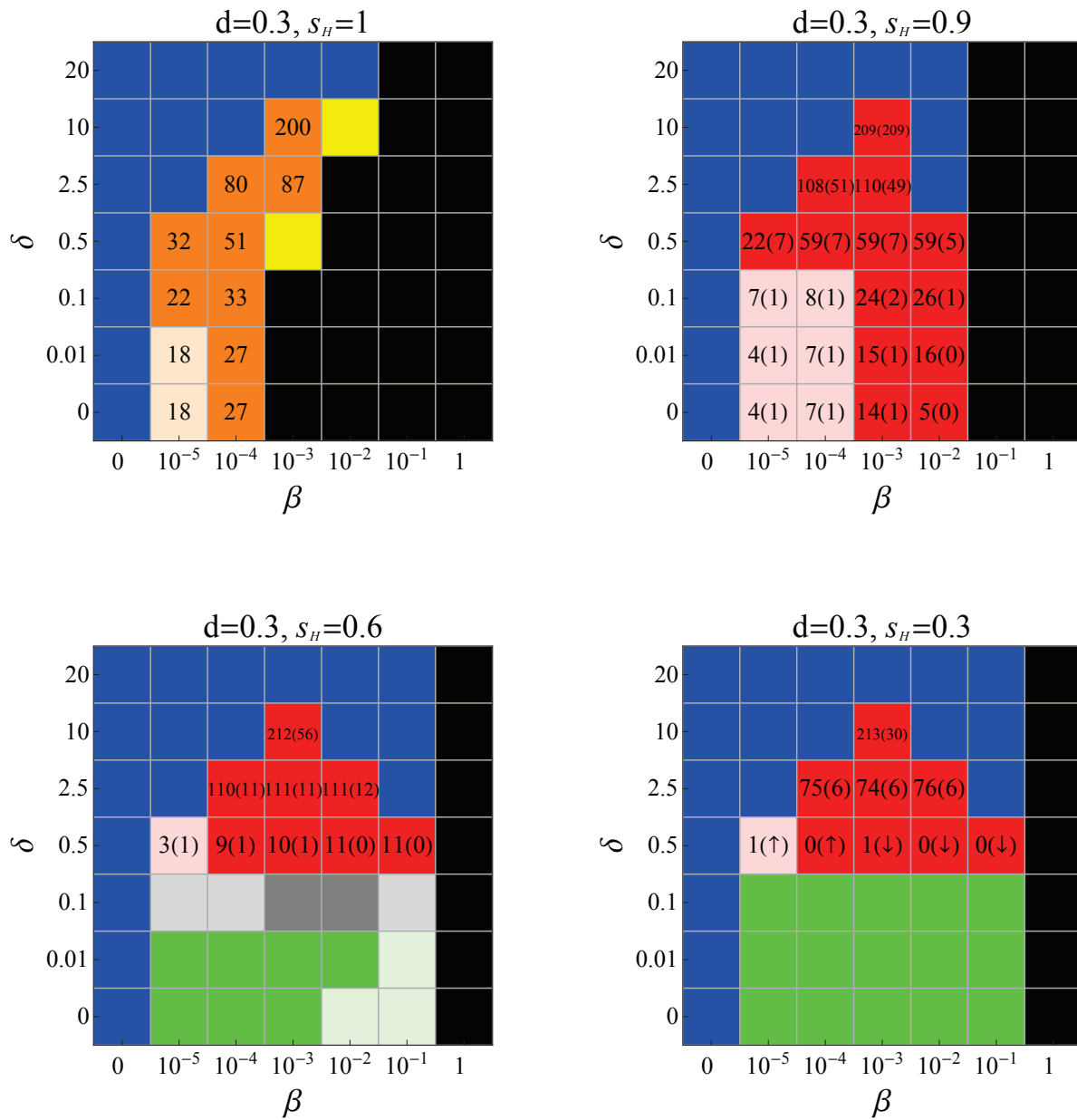




Figure SI 5.2.2



1 *Description of the results*

2 For each panel of Figures SI 5.2.1 and SI 5.2.2 the system of differential equations (1) and (2) were  
3 numerically evaluated for  $b = 1$ ,  $c_{H_1} = c_{P_2} = 0.05$ ,  $c_{H_2} = c_{P_1} = 0$  and the specified parameters  
4 over 500 time steps. The various possible coevolutionary scenarios of hosts and parasites are color-  
5 coded and summarized in detail in Tables SI 5.1.1 and SI 5.2.1; the latter table summarizes cases  
6 either slightly different for or even exclusive to the GFG model due to the additional possibility of  
7 infection. The squares featuring numbers represent the parameter combinations for which cycling  
8 around or immediate attainment of the fixed point of  $N^W$  is obtained. For  $s = 1$ , the amplitude  
9 of the cycles around the fixed point of  $N^W$  remains constant over time. For other values of  $s$  we  
10 distinguish two cases: We first evaluated the total amount of cycles around the fixed point of  $N^W$   
11 based on a numerical accuracy of 20 digits. Since the cycling dumps off in these cases, we also note  
12 in brackets the number of cycles until the amplitude is enclosed by the interval  $[0.95 N^W, 1.05 N^W]$ .

## 1 SI 6. Parametric plots for a matching-allele example

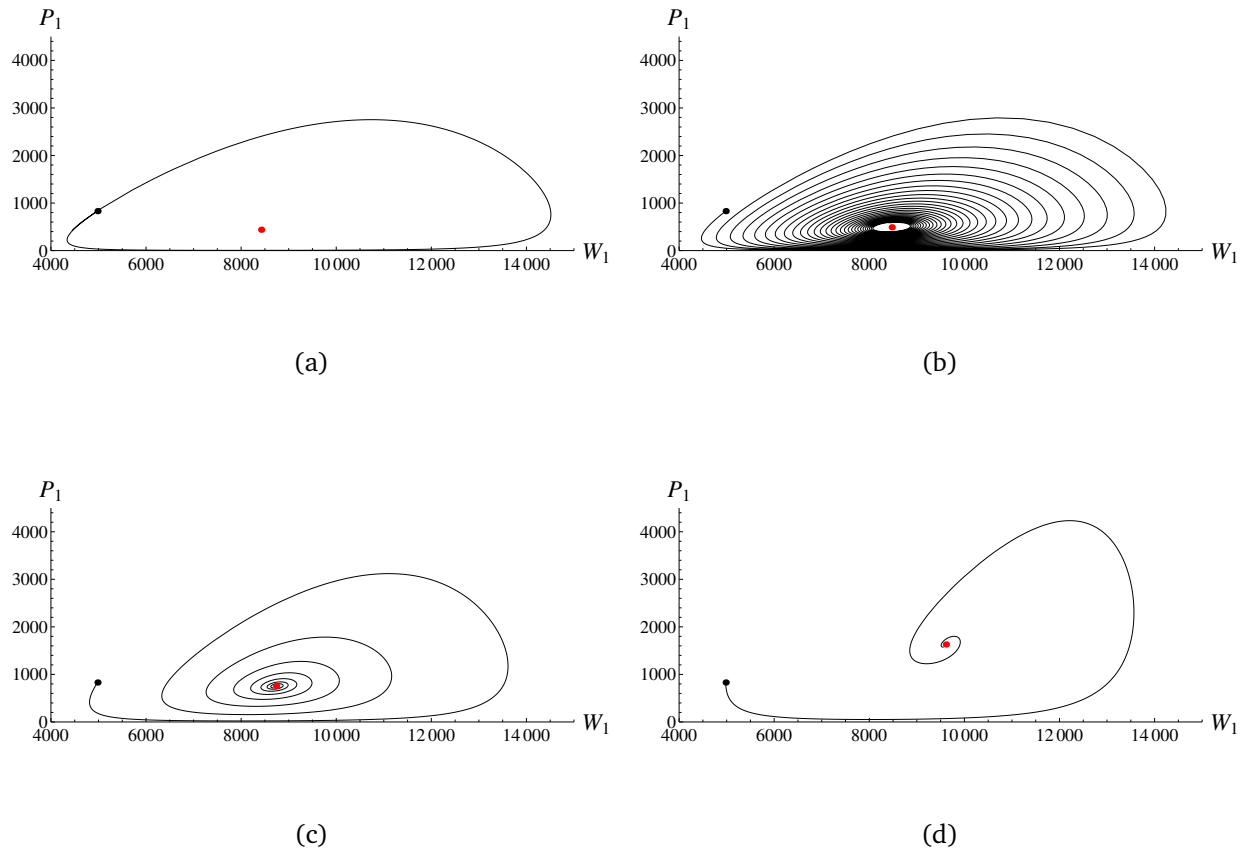
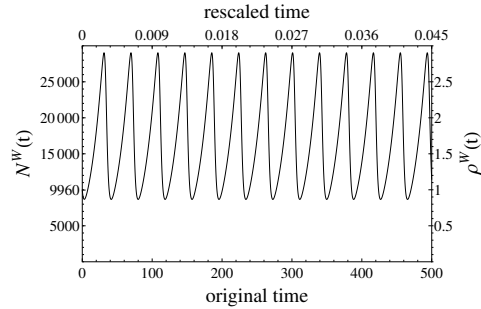
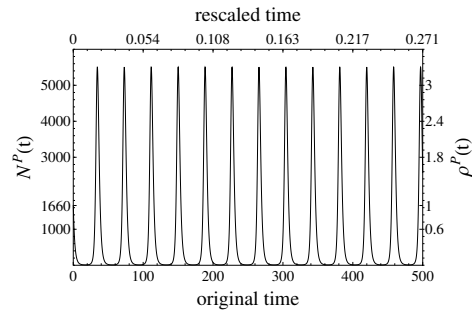


Figure SI 6.1 Parasite and host alleles of genotype one are plotted against each other for the MA model over time by numerically solving (1) and (2) for the following parameters (being equivalent for both genotypes):  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $\beta = 0.00012$  and  $c_P = c_H = 0.05$ . The initial conditions are  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$ . The selection coefficients  $s$  are given by (a) 1, (b) 0.9, (c) 0.6 and (d) 0.3. The parametric plots are shown for (a) 40, (b) 900, (c) 225 and (d) 76 time steps, which are the minimum amounts of time to complete one circle (a), come close to the fixed point (b), or to reach the fixed point (c), (d). The initial and fixed points are respectively colored in black and red.

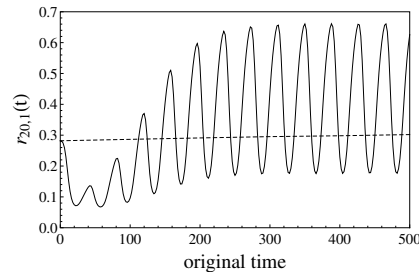
1 **SI 7. Time scaling and the distribution of singletons over time for a matching-allele**  
 2 **example**



(a)



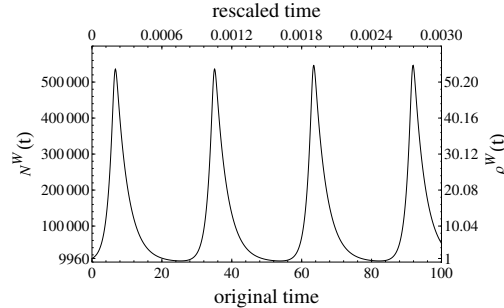
(b)



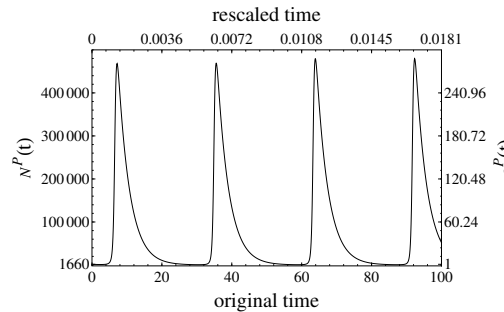
(c)

Figure SI 7.1 Population size changes in the host (a) and in the parasite (b) are generated for the MA model via the parameters  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $\beta = 0.00012$ ,  $c_P = c_H = 0.05$ , and  $s = 1$ . The initial conditions are  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$ , so that the reference population sizes  $N_{\text{ref}}$  of the host and the parasite are respectively given by 9960 and 1660 before their interaction starts at time zero. In both cases, the lower x-axes show time at the original scale of the dynamical system, whereas time is scaled by the respective values of  $N_{\text{ref}} \cdot 1/d$  for the upper x-axes. The left y-axes denote the absolute values of the changing population size and the right y-axes denote the relative values as  $\rho(t) = N(t)/N_{\text{ref}}$ . The number of infected alleles fluctuate between about ten and 5.500 individuals. We evaluated allelic spectra for every second generation (on the original scale) in both cases and (c) plot the relative number of singletons  $r_{20,1}(t)$  against time for the host (dashed) and the parasite (solid).

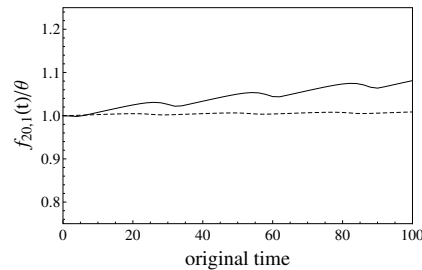
1 **SI 8. Time scaling and the distribution of singletons over time for a gene-for-gene**  
 2 **example**



(a)



(b)



(c)

Figure SI 8.1 Population size changes in the host (a) and in the parasite (b) are generated for the GFG model via the parameters  $b = 1$ ,  $d = 0.3$ ,  $\delta = 0$ ,  $\beta = 0.00001$ ,  $c_{H_1} = c_{P_2} = 0.05$ ,  $c_{H_2} = c_{P_1} = 0$  and  $s = 1$ . The initial conditions are  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$ , so that the reference population sizes  $N_{\text{ref}}$  of the host and the parasite are given by 9960 and 1660, respectively, before their interaction starts at time zero. In both cases, the lower x-axes show time at the original scale of the dynamical system, whereas time is scaled by the respective values of  $N_{\text{ref}} \cdot 1/d$  for the upper x-axes. The left y-axes denote the absolute values of the changing population size and the right y-axes denote the relative values as  $\rho(t) = N(t)/N_{\text{ref}}$ . The number of infected alleles fluctuate between about 720 and 478.000 individuals. We evaluated allelic spectra for every second generation (on the original scale) and (c) plot the absolute number of singletons scaled by  $\theta$ ,  $f_{20,1}(t)/\theta$ , against time for host (dashed) and parasite (solid).

## 1 SI 9. Distributions of singletons over time for various initial conditions

### 2 SI 9.1. A gene-for-gene example for 2%, 20% and 50% initially infected alleles

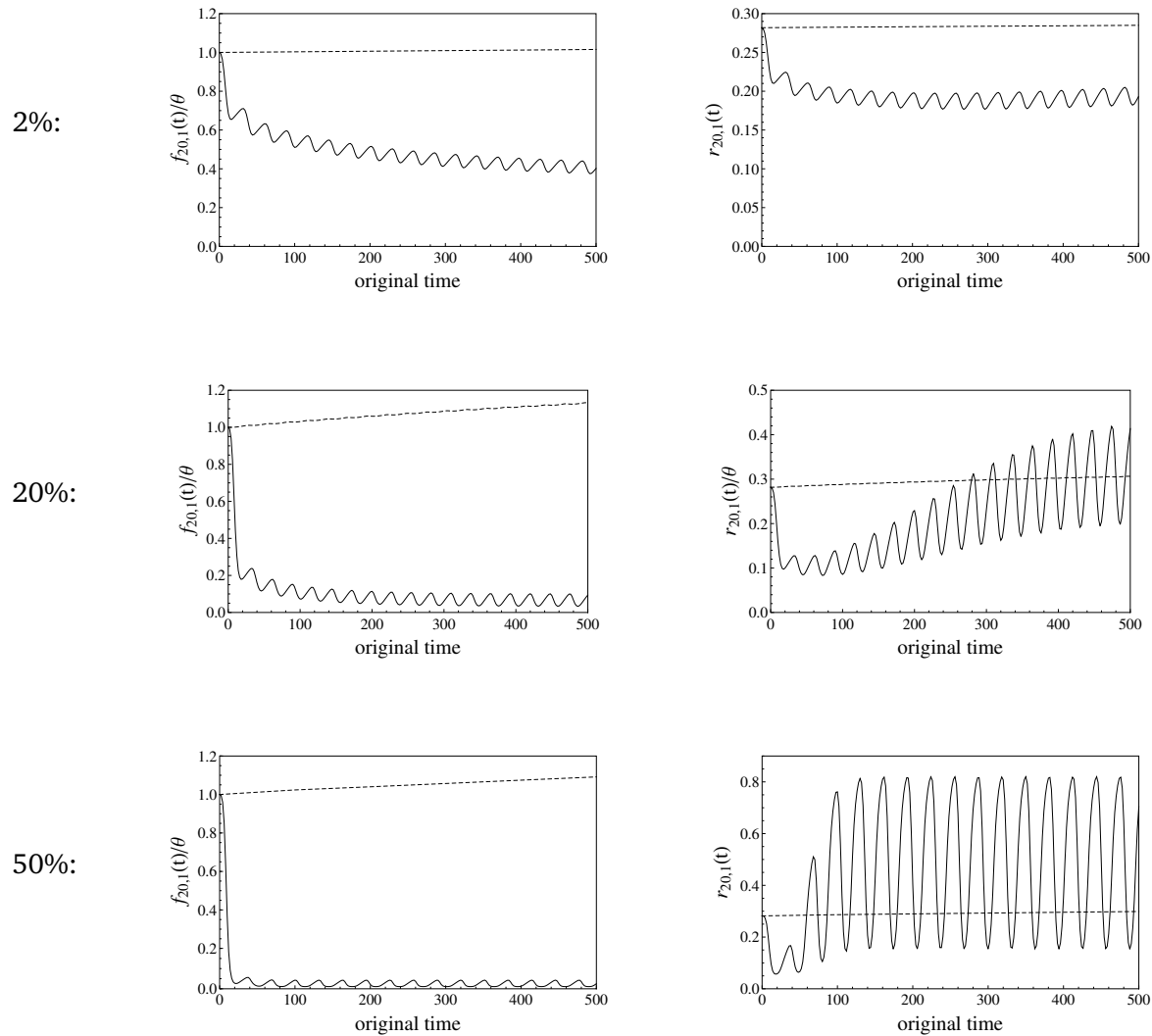


Figure SI 9.1 Population size changes in the host and in the parasite are generated for the GFG model via the parameters  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $c_{H_1} = c_{P_2} = 0.05$ ,  $c_{H_2} = c_{P_1} = 0$ ,  $s = 1$  and  $\beta = 0.000005$  (2%) or  $\beta = 0.00005$  (20% and 50%). The initial conditions are  $H_1 = H_2 = 49500$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 500$  (2%);  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$  (20%);  $H_1 = H_2 = 3333$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 833$  (50%). Therefore, the reference population sizes  $N_{\text{ref}}$  of the host and the parasite are, respectively, given by 101000 and 2000 (2%); 9960 and 1660 (20%); 9998 and 3332 (50%) before their interaction starts at time zero. We evaluated allelic spectra for every second generation (on the original scale) in all cases and plot the absolute number of singletons scaled by  $\theta$ ,  $f_{20,1}(t)/\theta$ , (left panels) and the relative number of singletons  $r_{20,1}(t)$  (right panels) against time for the host (dashed) and the parasite (solid).

1 SI 9.2. A matching-allele example for 2%, 20% and 50% initially infected alleles

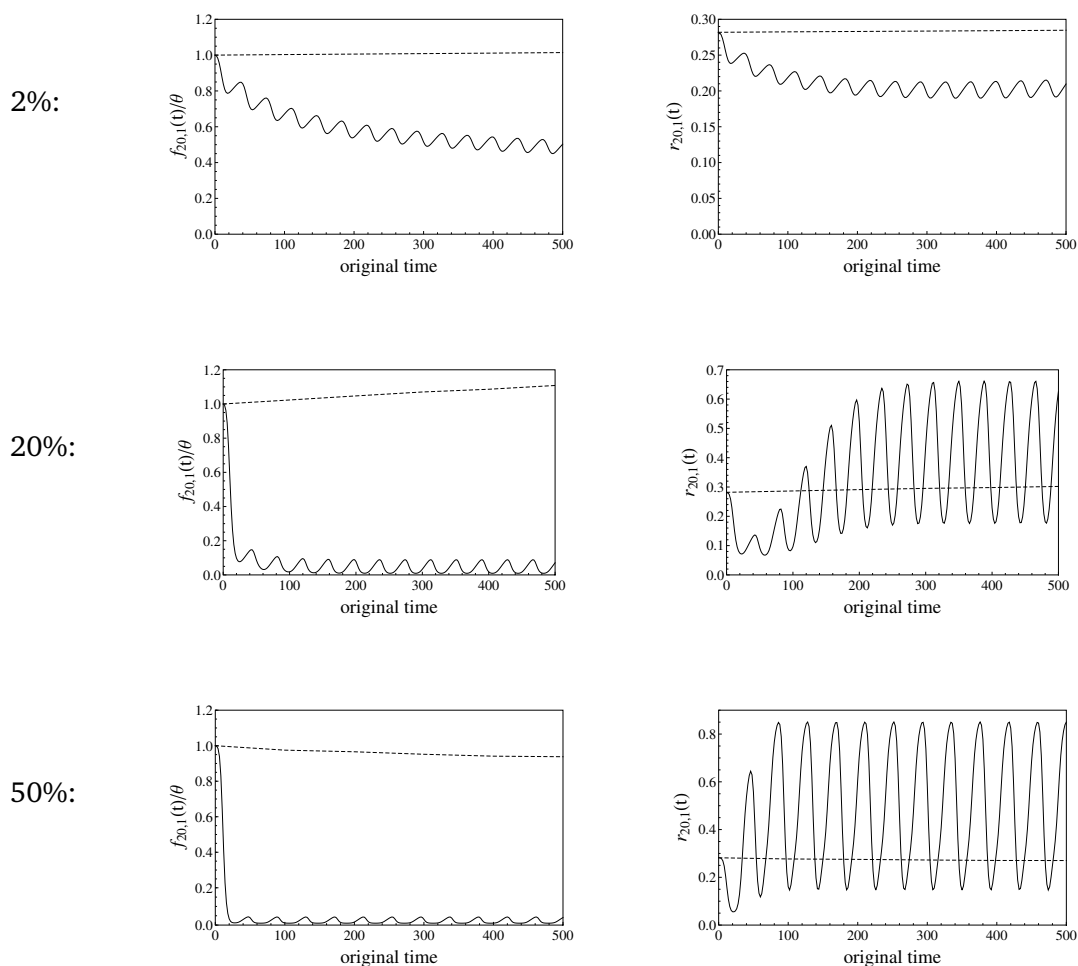


Figure SI 9.2 Population size changes in the host and in the parasite are generated for the MA model via the parameters  $b = 1$ ,  $d = 0.9$ ,  $\delta = 0.01$ ,  $c_H = c_P = 0.05$ ,  $s = 1$  and  $\beta = 0.0000012$  (2%),  $\beta = 0.000012$  (20%) and  $\beta = 0.00005$  (50%). The initial conditions are  $H_1 = H_2 = 49500$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 500$  (2%);  $H_1 = H_2 = 4150$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 415$  (20%);  $H_1 = H_2 = 3333$  and  $I_{11} = I_{12} = I_{21} = I_{22} = 833$  (50%). Therefore, the reference population sizes  $N_{\text{ref}}$  of the host and the parasite are, respectively, given by 101000 and 2000 (2%); 9960 and 1660 (20%); 9998 and 3332 (50%) before their interaction starts at time zero. We evaluated allelic spectra for every second generation (on the original scale) in all cases and plot the absolute number of singletons scaled by  $\theta$ ,  $f_{20,1}(t)/\theta$ , (left panels) and the relative number of singletons  $r_{20,1}(t)$  (right panels) against time for the host (dashed) and the parasite (solid).

1 **SI 10. The impact of multiple parasites per host and polycyclic diseases on detecting**  
2 **cycling population sizes for a gene-for-gene interaction**

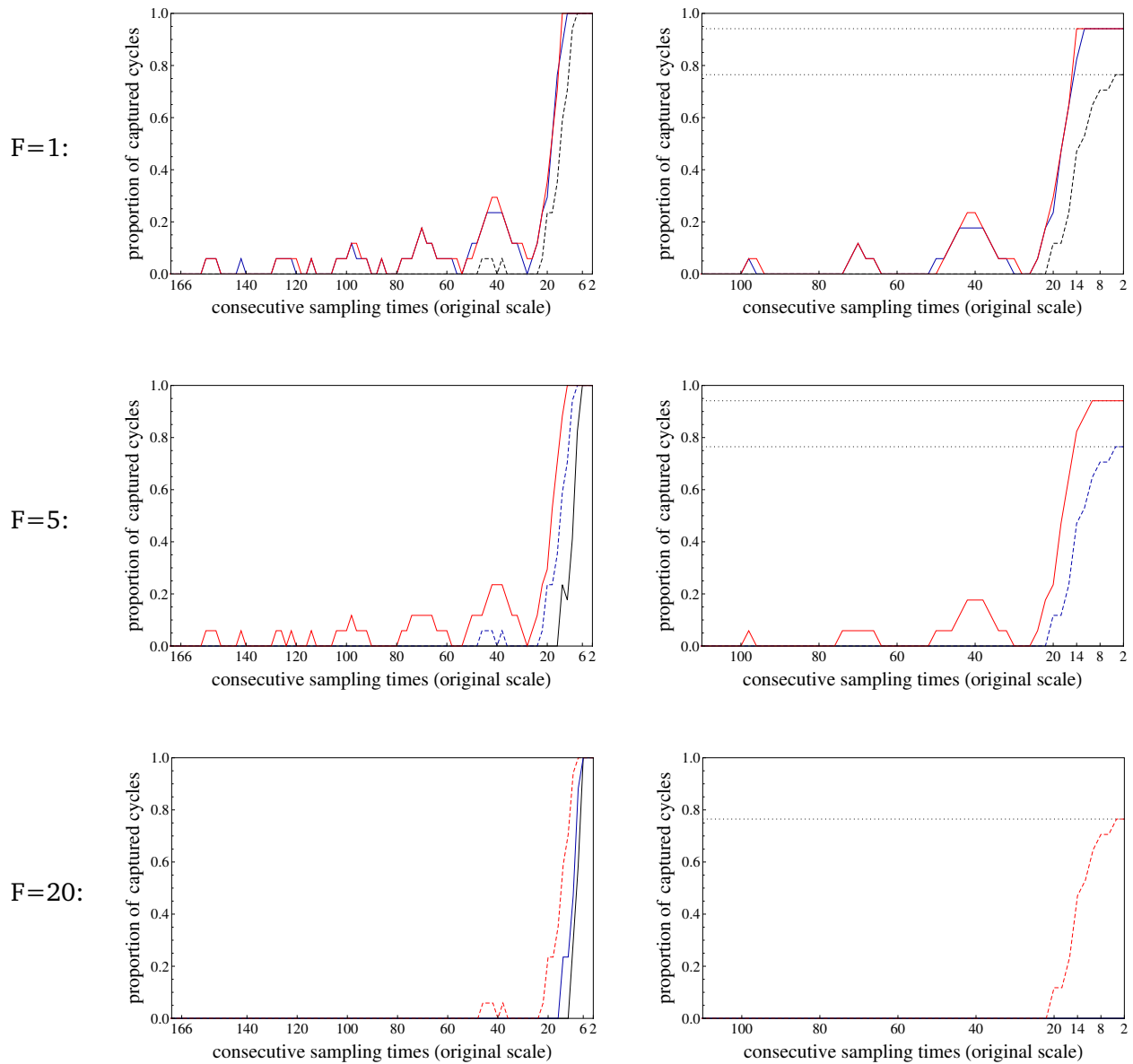


Figure SI 10.1



## 1 *Description of the results*

2 The dynamical system (1) and (2) was evaluated for the GFG model and the same parameter  
3 set as in the main example of *Results* (see Figure 3). Various values were considered for the number  
4 of parasite generations per host generation ( $E = 1$ , black curves;  $E = 5$ , blue curves;  $E = 20$ , red  
5 curves) and for the number of parasites per host ( $F$ ) taking the same values as  $E$ . Trajectories of  
6 the population size changes of the parasite were obtained for these nine parameter combinations  
7 and respectively employed into the analytical equation of the SFS, which was evaluated at every  
8 second generation (on the original scale) over an interval of 500 time points. Based on these  
9 datasets, we evaluated  $\Pi_{20}(j \cdot k/N_{\text{ref}})/\theta$  (see SI 4) for  $j = 2, 4, 6 \dots, 166$ , and  $k = 0, 1, 2, \dots$ , up to  
10  $j \cdot k$  taking the largest value equal or smaller than 500. Four time points are at least required to  
11 capture a single cycle and the value  $j = 166$ , for instance, means that a sample is taken every 166<sup>th</sup>  
12 generation at times zero, 166, 332 and 498. For every  $j_0$ , the times of alternating local minima and  
13 maxima of  $\Pi_{20}(j_0 \cdot k/N_{\text{ref}})$  were obtained over the values of  $k$ , before the rate of change as defined  
14 by  $\text{roc} = (\Pi_{20}((m+1)/N_{\text{ref}}) - \Pi_{20}(m/N_{\text{ref}}))/\Pi_{20}(m/N_{\text{ref}})$ ,  $m = 0, 1, 2, \dots$ , was applied until the end  
15 of the interval is reached. In all plots, the number of captured cycles are plotted (relative to the  
16 total number of cycles of the time interval) against consecutive sampling times  $j_0$ . For the left-hand  
17 panels, all cycles are considered, whereas for the right-hand panels only cycles are considered, if  
18 the rate of change between two consecutive extrema deviates by at least two percent, *i.e.*  $\text{roc} \geq 0.02$ .  
19 Every cycle is captured when sampling at every sixth generation on the left-hand side whereas not  
20 all cycles can be captured in any case and not at all for all chosen values of  $F < E$  on the right-  
21 hand side. The dashed lines illustrate cases with  $E = F$  giving equivalent curves. Also note that  
22 the number of captured cycles does not increase monotonically with the number of sampling times  
23 until sampling at every 20<sup>th</sup> generation in all the examples because a smaller number of sampling  
24 points (*e.g.* sampling at every 42<sup>nd</sup> generation in the first panel) can be more suitably distributed  
25 among the peaks and valleys of the considered time interval than a greater choice (*e.g.* sampling  
26 at every 30<sup>th</sup> generation in the first panel).