

# Needs a title

Sang Woo Park, Benjamin M Bolker

August 22, 2018

## 1 Introduction

The evolution of sexual reproduction presents a continuing question (Otto, 2009). Despite being the dominant mode of reproduction (Vrijenhoek, 1998), [BMB: *among?*] sexual reproduction entails numerous costs (Lehtonen et al., 2012). The most commonly mentioned is the cost of producing males (Smith, 1978). As males cannot produce offspring, sexual lineages is expected to be outgrown by their asexual counterpart that can grow as twice as fast (assuming that the sexual population produces 50% male and 50% female). This infamous *two-fold cost of sex* (Smith, 1978) relies on the assumption that everything else is equal. Then, what else is not equal and drives the sex to persist?

One explanation for the persistence of sexual reproduction is the Red Queen Hypothesis (Bell, 1982). The Red Queen Hypothesis suggests that sexually reproducing hosts overcome the cost of sex under strong parasite selection by producing genetically diverse offspring that are resistant to infection (Haldane, 1949; Jaenike, 1978; Hamilton, 1980). Host-parasite coevolution constantly creates selective advantage for rare genotypes, creating oscillation in genotypic frequencies, and allows for sexual reproduction to persist in the host population (Clarke, 1976; Hamilton, 1980).

Much of the theoretical work has focused on determining conditions under which parasite selection can maintain sexual reproduction in the host population. May and Anderson (1983) first noted that parasites must be extremely virulent to maintain sexual reproduction but later studies showed that sexual and asexual hosts can coexist even at intermediate virulence (Howard et al., 1994). Agrawal and Lively (2002) compared a wide range of infection genetics that determine parasite resistance and the dynamics that arise from different genetic architecture. Ashby and King (2015) showed that host genetic diversity also plays an important role in determining the strength of selection for sexual reproduction.

Some theoretical studies have departed from the classical population genetics framework to study effects of ecological and epidemiological structures on the Red Queen dynamics. Number of studies showed that incorporating ecological and epidemiological details can assist in supporting sexual reproduction in the host population (Galvani et al., 2001, 2003; Lively, 2009, 2010b). In contrary to

these findings, MacPherson and Otto (2017) showed that Red Queen dynamics (i.e., cycles in allele frequencies) fail to persist when explicit epidemiological structure is taken into account with coevolutionary dynamics.

On the other hand, empirical studies have mostly focused on confirming predictions that stem from the Red Queen Hypothesis. Typical among them are local adaptation, time-lagged selection, and association between parasite prevalence and host reproductive mode (see Tobler and Schlupp (2008) and Vergara et al. (2014b) for reviews). A key example is the snail population in New Zealand that serve as an intermediate hosts for trematodes [CITE]. Through several decades of work, Lively *et al.* demonstrated that the population satisfies necessary conditions for the host-parasite coevolutionary dynamics and provide support for the hypothesis (Lively, 1987, 1989; Dybdahl and Lively, 1995, 1998; Jokela et al., 2009; Vergara et al., 2014b; Gibson et al., 2016). While many studies provide only indirect evidence, recent studies show that more direct evidence can be achieved using experimental systems (Auld et al., 2016; Slowinski et al., 2016).

Even though the Red Queen Hypothesis has gained some support both theoretically and empirically, there still remains a gap between theory and data. Many models for Red Queen Hypothesis rely on simplifying assumptions that are not applicable to natural populations and make predictions based on assumed parameters. In particular, none of the Red Queen models reviewed by Ashby and King (2015) use statistical tools to relate model to data. However, there are exceptions: Lively (1992) postulated that infection prevalence should be positively correlated with frequency of sexual hosts and later formalized the idea with a mathematical model (Lively, 2001). The prediction has since been confirmed by many empirical studies, most of which are based on the snail-trematode system (Lively and Jokela, 2002; Kumpulainen et al., 2004; Vergara et al., 2013; McKone et al., 2016; Gibson et al., 2016). Surprisingly such correlation was not observed in a different snail-trematode system (Dagan et al., 2013a). [BMB: *A little more on qualitative requirements?*] [SWP: *What do you mean?*]

Here, we try to bridge the gap between theory and data further. We extend the model used by Lively (2010b) to account for demographic stochasticity and simple population structure. Then, we fit the model to observational data from Dagan et al. (2013a); McKone et al. (2016); Vergara et al. (2014b) using Approximate Bayesian Computation (ABC) to estimate biologically relevant parameters. Using estimated parameters, we assess model fits and perform a power analysis to test the prediction that infection prevalence is positively correlated with frequency of sexual reproduction Lively (2001). [BMB: *More on power analysis*]

## 2 Methods

### 2.1 Data

We consider observational data from two snail-trematode populations in New Zealand (Vergara et al., 2014b; McKone et al., 2016) and a similar snail-trematode population in Israel (Dagan et al., 2013a). The snail-trematode system has been extensively studied under the context of the Red Queen Hypothesis so we expect a simple Red Queen model to fit reasonably well. Data collected by Dagan et al. (2013a) and Vergara et al. (2014b) were obtained from their Dryad repositories (Dagan et al., 2013b; Vergara et al., 2014a) and data collected by McKone et al. (2016) was extracted from their figure.

### 2.2 Model

**[SWP: TODO: read Lively 2018]** We model obligately sexual hosts competing with obligately asexual hosts in a meta-population by extending the model introduced by Lively (2010b). The model is a discrete time susceptible-infected (SI) model with natural mortality and virulence (defined as reduction in offspring production among infected hosts). It is a suitable candidate model for this study as it captures essential structures that are present in basic epidemiological and population genetics models and is general enough to be applied to broad range of natural systems. We do not consider mechanistic details of the snail-trematode system such as life history of trematodes (Vergara et al., 2014b). We incorporate population structure and allow for mixing between populations. Each population can be equivalently considered as a sampling site in the observed populations.

All hosts are assumed to be diploids with two biallelic loci, and parasites are assumed to be haploids. Let  $S_{ij}^k(t)$  and  $A_{ij}^k(t)$  be the number of sexual and asexual hosts with genotype  $ij$  from population  $k$  at generation  $t$ . For simplicity, we drop the superscript representing population and write  $S_{ij}(t)$  and  $A_{ij}(t)$ ; every population is governed by the same set of equations unless noted otherwise (e.g., when we account for interaction between populations). Following Lively (2010b), the expected amount of genotypic contribution (before recombination or outcrossing) by sexual hosts is given by

$$S'_{ij} = c_b(1 - s)(W_U S_{ij,U}(t) + W_I S_{ij,I}(t)), \quad (1)$$

where  $s$  is the proportion of males produced by sexual hosts, and  $S_{ij,U}$  and  $S_{ij,I}$  are the number of uninfected and infected sexual hosts in a population.  $W_U$  and  $W_I$  represent their corresponding fitnesses where virulence is defined as  $V = 1 - W_I/W_U$ . We allow for cost of sex to vary by multiplying a scale parameter,  $c_b$ , to the growth term, where  $2/c_b$  corresponds to a two fold cost of sex (Ashby and King, 2015). Recombination and outcrossing are modeled after incorporating genotypic contributions from other populations.

We define

$$W_U = \frac{b_U}{1 + a_U N(t)}, W_I = \frac{b_I}{1 + a_I N(t)}$$

where  $b_U$  and  $b_I$  are number of offspring produced by uninfected and infected hosts, respectively, and  $a_U$  and  $a_I$  determine their corresponding density dependent effects (Lively, 2010b; Smith and Slatkin, 1973). For simplicity, we assume that  $a_U = a_I$  so that virulence can be defined strictly in terms of decrease in offspring production and is constant for any density:  $V = 1 - b_I/b_U$ . [BMB: more specific] [SWP: Lively doesn't say much beyond this... I think it's OK?]

Asexual hosts are assumed to be strictly clonal. Then, the expected amount of genotypic contribution by asexual hosts is given by

$$A'_{ij} = W_U A_{ij,U}(t) + W_I A_{ij,I}(t), \quad (2)$$

where  $A_{ij,U}$  and  $A_{ij,I}$  are the number of uninfected and infected asexual hosts in a population.

We assume that proportion  $\epsilon_{\text{mix}}$  of a population mix with other populations. Then, the expected number of offspring in the next generation (accounting for contributions from all populations) is given by

$$\begin{aligned} E(S_{ij}^k(t+1)) &= f_{\text{sex}} \left( (1 - \epsilon_{\text{mix}}) (S_{ij}^k)' + \frac{\epsilon_{\text{mix}}}{n_{\text{pop}} - 1} \sum_{h \neq k} (S_{ij}^h)' \right), \\ E(A_{ij}^k(t+1)) &= (1 - \epsilon_{\text{mix}}) (A_{ij}^k)' + \frac{\epsilon_{\text{mix}}}{n_{\text{pop}} - 1} \sum_{h \neq k} (A_{ij}^h)', \end{aligned} \quad (3)$$

where  $f_{\text{sex}}(x)$  is the function that models sexual reproduction, including recombination probability  $r_{\text{host}}$  and outcrossing, and  $n_{\text{pop}}$  is the number of populations modeled. Then, the total number of sexual and asexual hosts in the next generation given by Poisson random variables with mean specified previously. We also allow for stochastic migration to avoid fixation:

$$\begin{aligned} S_{ij}^k(t+1) &\sim \text{Poisson}(\lambda = E(S_{ij}^k(t+1))) + \text{Bernoulli}(p = p_{ij,\text{sex}}), \\ A_{ij}^k(t+1) &\sim \text{Poisson}(\lambda = E(A_{ij}^k(t+1))) + \text{Bernoulli}(p = p_{ij,\text{asex}}), \end{aligned} \quad (4)$$

where  $p_{ij,\text{sex}}$  and  $p_{ij,\text{asex}}$  are the probabilities of a sexual and an asexual host with genotype  $ij$  entering a population.

[BMB: No epistasis?] [SWP: Is this clearer?] Infection is modeled using the matching alleles model (Otto and Michalakis, 1998). We assume that snails are equally susceptible to parasites that match either haplotype. However, parasites must carry same alleles in both loci in order to match a host haplotype. The total number of infected hosts that carry parasite with genotype  $i$  at generation  $t$  is given by:

$$I_i(t) = \sum_p 2^{\delta_{ij}} (S_{ij,i,I}(t) + A_{ij,i,I}(t)), \quad (5)$$

where a  $\delta_{ij}$  is Kronecker delta.  $\delta_{ij}$  equals 1 when  $i = j$  and 0 otherwise.  $S_{ij,i,I}(t)$  and  $A_{ij,i,I}(t)$  represent the expected numbers of sexual and asexual hosts that have genotype  $ij$  and are infected with genotype  $i$  parasite. Following Ashby

and King (2015), we assume that mutation can occur in one locus with probability  $r_{\text{parasite}}$ . Mutation is modeled using a deterministic process, as we introduce stochasticity during the actual infection process. We also allow for stochastic external migration of an infected host carrying parasite  $i$  with probability  $p_{i,\text{parasite}}$  to avoid fixation.

The total expected number of infectious contacts made by infected hosts within a population is given by  $\lambda_i^k = \beta^k I_i^k(t)$ , where  $\beta^k$  is the transmission rate of each population, and  $I_i^k(t)$  is the number of infected hosts accounting for mutation and migration. Since we allow for mixing between populations, infected hosts can make contact with susceptible hosts in other populations. **[SWP: TODO: explain mixing]** Then, the total amount of infectious contacts, coming from hosts that carry genotype  $i$  parasite, that is received by susceptible hosts in population  $k$  is given by

$$\lambda_{i,\text{total}}^k = (1 - \epsilon_{\text{mix}})\lambda_i^k + \frac{\epsilon_{\text{mix}}}{n_{\text{pop}} - 1} \sum_{l \neq k} \lambda_i^l \quad (6)$$

Then, the force of infection that a susceptible host with genotype  $ij$  experiences in generation  $t + 1$  is given by

$$\text{FOI}_{ij}^k = \frac{\lambda_{i,\text{total}}^k + \lambda_{j,\text{total}}^k}{2N^k(t + 1)}, \quad (7)$$

where  $N^k(t + 1) = \sum_{i,j} S_{ij}^k(t + 1) + A_{ij}^k(t + 1)$  is the total number of hosts in generation  $t + 1$ . The probability that a susceptible host with genotype  $ij$  in population  $k$  becomes infected in the next generation is given by

$$P_{ij}^k(t + 1) = 1 - \exp\left(-\text{FOI}_{ij}^k\right). \quad (8)$$

Finally, number of infected hosts in the next generation is determined by a binomial process:

$$\begin{aligned} S_{ij,I}^k(t + 1) &\sim \text{Binom}(S_{ij}^k(t + 1), P_{ij}^k(t + 1)), \\ A_{ij,I}^k(t + 1) &\sim \text{Binom}(A_{ij}^k(t + 1), P_{ij}^k(t + 1)). \end{aligned} \quad (9)$$

The expected number of infected hosts that have genotype  $ij$  and are infected by parasites with genotype  $i$  in the next generation is given by a ratio of  $\lambda$ :

$$\begin{aligned} S_{ij,i,I}^k(t + 1) &= \frac{\lambda_{i,\text{total}}^k}{\lambda_{i,\text{total}}^k + \lambda_{j,\text{total}}^k} S_{ij,I}^k(t + 1) \\ A_{ij,i,I}^k(t + 1) &= \frac{\lambda_{i,\text{total}}^k}{\lambda_{i,\text{total}}^k + \lambda_{j,\text{total}}^k} A_{ij,I}^k(t + 1) \end{aligned} \quad (10)$$

### 2.3 Simulation design and parameterization

Many Red Queen models have focused on competition between a single asexual genotype and multiple sexual genotypes or have assumed equal genetic diversity

between asexual and sexual hosts (see (Ashby and King, 2015) for a review of previous Red Queen models) but neither of these assumptions are realistic. Instead, Ashby and King (2015) adopted a more realistic approach by allowing for stochastic migration of an asexual genotype to a population. Here, we combine these methods. We allow for stochastic external migration of asexual hosts with different genotypes into the system but fix the number of asexual genotypes (denoted by  $G_{\text{asex}}$ ) that can be present in the system. The number of sexual genotypes ( $G_{\text{sex}}$ ) that can be present in the population remains equal to the size of the genotypic space ( $= 10$  for diploid hosts with two biallelic loci).

[BMB: *Explain?*] Given a value for  $G_{\text{asex}}$ , asexual genotypes that can be introduced to the population are uniformly chosen from the entire genotypic space in the beginning of the simulation. Limiting asexual genotypes account for difference in genetic diversity between asexual and sexual lineages. We estimate  $G_{\text{asex}}$  to test whether greater asexual genetic diversity can be supported.

To account for differing number of sexual and asexual genotypes, we let

$$p_{ij,\text{sex}} = 1 - (1 - p_{\text{host}})^{1/G_{\text{sex}}},$$

$$p_{ij,\text{asex}} = \begin{cases} 1 - (1 - p_{\text{host}})^{1/G_{\text{asex}}} & \text{if } ij \in \{\text{asexual genotypes}\} \\ 0 & \text{otherwise} \end{cases}, \quad (11)$$

where  $p_{\text{host}}$  is the probability that at least one sexual and asexual host enters the population in a generation. We scale the probability of infected host carrying parasite genotype  $i$  in a similar way for interpretability:

$$p_{i,\text{parasite}} = 1 - (1 - p_{\text{infected}})^{1/4}, \quad (12)$$

where  $p_{\text{infected}}$  is the probability that at least one infected host enters the population in a generation.

Each simulation consists of 40 populations. Every population is initialized with 2000 sexual hosts where 80 of them are infected. They are assumed to be in Hardy-Weinberg equilibrium with ratio between alleles being exactly half. Transmission rate,  $\beta^k$ , is randomly drawn for each population from a gamma distribution with mean  $\beta_{\text{mean}}$  and coefficient of variation  $\beta_{\text{cv}}$ . Simulation runs for 500 generations without introduction of asexuals. At generation 501, 10 asexual hosts of a single genotype are introduced to each population (note that asexual genotype introduced can vary across population) and simulation runs for 600 generations while allowing for stochastic migration of asexuals.

## 2.4 Approximate Bayesian Computation

[BMB: *More on probe matching; Kendall et al.?*] [SWP: *Not clear what you're looking for...*]

We use Approximate Bayesian Computation (ABC) to fit the model (Toni et al., 2009). ABC relies on comparing summary statistics of observed data and those of simulated data and is particularly useful when the exact likelihood function is not available. We consider mean proportion of infected and sexually

reproducing snails in the system and variation in these proportions – measured by coefficient of variation (CV) – across space (population) and time as our focal summary statistics. These summary statistics are calculated for both observed and simulated data and are used in ABC. As Dagan et al. (2013a) and McKone et al. (2016) only reported proportion of males, proportion of sexual hosts are assumed to be twice proportion of males.

CV across space is calculated by first calculating mean proportions by averaging across time (generation) for each site (population) and then taking the CV of these mean proportions. CV across time (generation) is calculated by first averaging proportions across space (population) at each generation and then taking the CV. For purely spatial data (Dagan et al. (2013a) and McKone et al. (2016)), CV across space is calculated without averaging across time. Sampling error is not taken into account when summary statistics are calculated from simulated populations.

We use weakly informative priors for all parameters that we estimate except  $c_b$ , a scale parameter for the cost of sex (see Table 1 for prior distributions used and parameters assumed). The prior distribution for the scale parameter is chosen so that 95% quantile of cost of sex ( $2/c_b$ ) is approximately equal to the 95% confidence interval reported by Gibson et al. (2017). All other parameters are assumed to be fixed for simplicity.

We start by performing basic ABC. For each random parameter sample drawn from the prior distribution, the model is simulated and a sample of simulated populations is drawn from the simulated system such that the number of sampled population is equal to the number of sites collected in a study. Then, summary statistics are calculated based on the last 100 generations out of 1100 generations and the parameter is accepted if the distance between simulated and observed data is less than a tolerance value. Distance is measured by the sum of absolute differences in summary statistics between simulated and observed data. This process is repeated until 100 parameter sets are accepted.

After the first run ( $t = 1$ ), equal weights ( $w_{i,1} = 1/100$ ) are assigned to each accepted parameter set  $\theta_{i,1}$ , where  $1 \leq i \leq 100$ . For any run  $t > 1$ , a weighted random sample ( $\theta^*$ ) is drawn from the accepted parameters of the previous run ( $t - 1$ ) with weights  $w_{i,t-1}$  and a parameter sample ( $\theta_{i,t}$ ) is proposed from a multivariate normal distribution with a mean  $\theta^*$  and a variance covariance matrix that is equal to  $\sigma_{t-1}^2 = 2\text{Var}(\theta_{1:N,t-1})$ , where  $\text{Var}(\theta_{1:N,t-1})$  is the weighted variance covariance matrix of the accepted parameters from the previous run.  $N$  is the total number of accepted parameters from the previous run.

$G_{\text{asex}}$  is rounded to the nearest integer and the model is simulated. If a proposed parameter is accepted, the following weight is assigned:

$$w_{i,t} = \frac{\pi(\theta_{i,t})}{\sum_{j=1}^{100} w_{j,t-1} q(\theta_{j,t-1} | \theta_{i,t}, \sigma_{t-1}^2)}$$

where  $\pi(\cdot)$  is a prior density and  $q(\cdot | \theta_{i,t}, \sigma_{t-1}^2)$  is a multivariate normal density with mean  $\theta_{i,t}$  and variance covariance matrix  $\sigma_{t-1}^2$ . For each run, 100 parameters are accepted and weights are normalized at the end to sum to 1.

| Notation                | Description   | Prior distribution/parameter values          | Source               |
|-------------------------|---|--|----------------------|
| $\beta_{\text{mean}}$   | Mean transmission rate  | Gamma( $k = 2, \theta = 10$ )                | Assumption           |
| $\beta_{\text{CV}}$     | CV transmission rate  | Gamma( $k = 2, \theta = 0.5$ )               | Assumption           |
| $V$                     | Virulence   | Beta( $\alpha = 6, \beta = 2$ )              | Assumption           |
| $\epsilon_{\text{mix}}$ | Mixing proportion   | Beta( $\alpha = 1, \beta = 9$ )              | Assumption           |
| $G_{\text{asex}} - 1$   | Number of asexual genotypes - 1   | BetaBinomial( $N = 9, p = 3/9, \theta = 5$ ) | Assumption           |
| $c_b$                   | Cost of sex scale   | LogNormal( $\mu = -0.07, \sigma = 0.09$ )    | Gibson et al. (2017) |
| $s$                     | Proportion of male offsprings produced                                      | 0.5  | Assumption           |
| $b_U$                   | Number of offsprings produced by an uninfected host                         | 20   | Lively (2010b)       |
| $b_I$                   | Number of offsprings produced by an infected host                           | $(1 - V)b_U$                                 | Lively (2010b)       |
| $a_U$                   | Density dependent effect coefficient of uninfected hosts                    | 0.001  | Lively (2010b)       |
| $a_I$                   | Density dependent effect coefficient of infected hosts                      | 0.001  | Lively (2010b)       |
| $r_{\text{host}}$       | Host recombination probability  | 0.2  | Lively (2010b)       |
| $r_{\text{parasite}}$   | Parasite mutation probability   | 0.05   | Assumption           |
| $p_{\text{host}}$       | Probability that at least one sexual and asexual host enters the population | 0.1  | Assumption           |
| $p_{\text{infected}}$   | Probability that at least one infected host enters the population           | 0.02   | Assumption           |

Table 1: **Parameter descriptions and values.** Parameters with prior distributions are estimated via Approximate Bayesian Computation (ABC).  $k$  and  $\theta$  in Gamma distribution represent shape and scale parameters where mean and squared CV are given by  $k\theta$  and  $1/k$ , respectively.  $\alpha$  and  $\beta$  in Beta distribution represent shape parameters where mean and squared CV are given by  $\alpha/(\alpha + \beta)$  and  $\beta/(\alpha^2 + \alpha\beta + \alpha)$ .  $N$ ,  $p$  and  $\theta$  in Beta binomial distributions represent number of trials, probability of success, and overdispersion parameters (Morris et al., 1983). We define prior on  $G_{\text{asex}} - 1$  instead to always maintain at least one asexual genotype in the system.  $\mu$  and  $\sigma$  in log-normal distribution represent mean and standard deviation on a log scale. All other parameters are fixed throughout simulations.



This method, known as the Population Monte Carlo approach (Turner and Van Zandt, 2012), allows for sampling more efficiently while ensuring that final result still satisfies criteria to be a correct (approximate) Bayesian posterior. All statistical results reported are weighted by parameter weights of the final run.

For each observed datum, we perform 4 runs with decreasing tolerance every run. For spatial data (Dagan et al., 2013a; McKone et al., 2016), four summary statistics are compared: mean proportion of infected and sexually reproducing snails and CV in these proportions across populations. Tolerance values of 1.6, 0.8, 0.6 and 0.4 are used for each run. For spatiotemporal data (Vergara et al., 2014b), six summary statistics are compared: mean proportion of infected and sexually reproducing snails, CV in these proportions across populations and CV in these proportions across generations. Larger tolerance values (2.4, 1.2, 0.9 and 0.6) are used for each run to account for higher number of summary statistics being compared. Tolerance value of the final run is chosen so that a parameter set will be accepted if its each simulated summary statistic deviates from the corresponding observed summary statistic by 0.1 units on average. First three tolerance values are chosen in a decreasing order to reach the final step quicker. *[BMB: Intuition for biological meaning of distances?] [SWP: I don't think there is one?] [BMB: What are the summary stats?] [SWP: They're explained in the beginning of the section]*

## 2.5 Power analysis

Using estimated parameters for each data, we calculate the power to detect a correlation between infection prevalence and frequency of sexual hosts. For each parameter sample from the final run of the ABC, 10 simulations are ran. For each simulation, we start by setting a reference generation to 1001st generation choosing  $n$  populations at random from 40 simulated populations. For each selected population, hosts are divided into four categories based on their infection status (infected/uninfected) and reproductive mode (asexual/sexual), and mean proportion of hosts in each category is calculated by averaging over two consecutive generations. We assume that a year contains of two snail generations [CITE] and that samples are taken within a short period of time.

Independent multinomial samples of size  $m$  are drawn from each selected population based on the proportions in each four categories. Correlation between proportion of infected hosts and proportion of sexual hosts is tested using Spearman's rank correlation at 5% significance level. We repeat the process 100 times by changing the reference generation from 1001st generation to 1099th generation.

## 3 Results

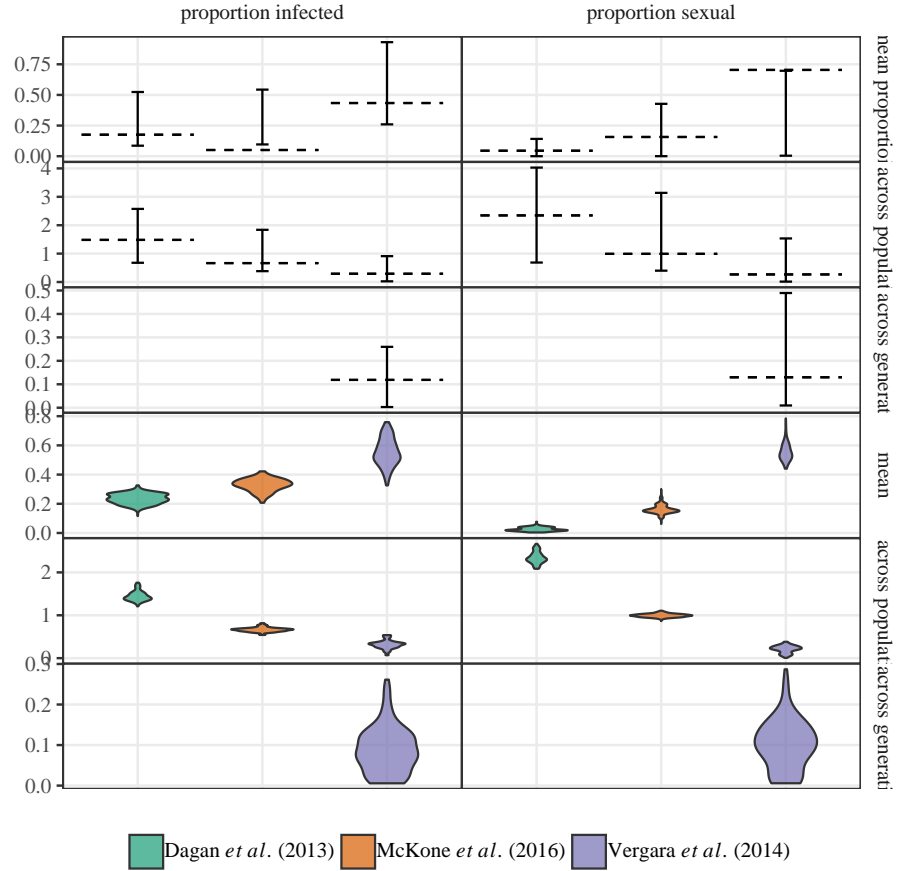
Fig. 1 compares observed summary statistics with fitted and predicted summary statistics. Fitted summary statistics are those that are accepted via ABC and

can be interpreted as underlying summary statistics of the study sites estimated by the model. [SWP: *Is this interpretation OK?*] Predicted summary statistics are obtained by simulating from estimated parameters and represent what could have been the underlying summary statistics if other sites were chosen from the system. As we account for uncertainty in unobserved sites by simulating a greater number of populations, there is large variation in predicted summary statistics.

We find that our simple meta-population Red Queen model can capture observed variation in infection prevalence and frequency of sexual hosts reasonably well; both temporal and spatial variation (measured by CV across mean proportions) are well-matched by the model. On the other hand, the model tends to overestimate mean proportion of infected hosts. Dagan et al. (2013a), McKone et al. (2016) and Vergara et al. (2014b) reported mean infection prevalence of 17.5%, 5.1% and 44% in their study sites, respectively. The *estimated* mean (95% quantile) infection prevalence is 24.0% (17.4% - 28.6%), 31.3% (23.3% - 40.4%) and 54.2% (36.0% - 73.0%), respectively. The model also underestimates mean frequency of sexual hosts for Dagan et al. (2013a) and Vergara et al. (2014b) study sites. Observed mean frequency of sexual hosts is 4.5% and 70.4%, respectively, whereas corresponding estimated mean (95% quantile) are 2.6% (0.7% - 4.8%) and 59.6% (N.A - 67.9%). [SWP: *wquant doesn't give us a value at 2.5%. What should I do?*] As model fitting is performed by minimizing the sum of absolute distance between observed and simulated summary statistics, our method does not guarantee all summary statistics to be equally well fitted.

To further diagnose the fit, we compare the predicted relationships between mean infection prevalence and mean frequency of sexual hosts in each population (averaged over last 100 generations) with the observed data (Fig. 2). Note that Fig. 2 appears to be more variable than Fig. 1 as it plots density of all simulated populations and hence accounts for uncertainty in unsampled populations. Despite being able to reproduce the summary statistics reported by Dagan et al. (2013a) well, our model is unable to capture the qualitative trend between proportion of sexual hosts and proportion of infected hosts (Fig. 2; Dagan et al. (2013a)). Both simulated data and observed data mostly consist of asexual populations but our model predicts sexual reproduction to be maintained when infection prevalence is high ( $> 40\%$ ). On the other hand, Dagan et al. (2013a) data suggests that sexual reproduction is only maintained when infection prevalence is low ( $< 20\%$ ). Similarly, overestimation of infection prevalence is strongly pronounced in our prediction of system studied by McKone et al. (2016).

While there are a few data points that appear to be outliers compared to our predictions for Vergara et al. (2014b), it is important to note that Fig. 2 does not capture temporal variation as we average over 100 generations to obtain the “mean” relationship. The observed data are more likely to be samples across a few generations and the cyclic nature of the Red Queen dynamics is likely to have created more variation in the data. On the other hand, Vergara et al. (2014b) reported greater than 90% sexual snails throughout 5 years in one of



**Figure 1: Summary statistics of the observed data vs. distribution of summary statistics of the simulated data from the posterior samples.**

Dotted horizontal line represents observed summary statistics. Violin plots show weighted distribution of fitted summary statistics (i.e., summary statistics that were accepted during Approximate Bayesian Computation). Error bars show 95% weighted quantiles of predicted summary statistics. For each posterior sample, 10 simulations are run and each simulated system is sampled at random 100 times so that each sample consists of equal number of populations as number of sites in fitted data. Then, summary statistics are calculated for each sample and are weighted by their corresponding weights.

their study sites but it seems unusually high based on our model prediction.

We find that there is a region (around 30% infection prevalence) in which proportion of infected hosts remains almost constant while proportion of sexual hosts increases (most clearly visible in the fits to McKone et al. (2016) and

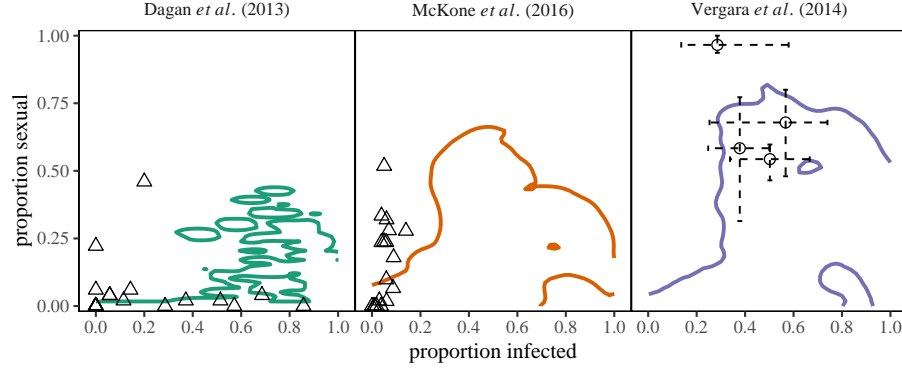


Figure 2: **Predicted relationship between mean infection prevalence and mean proportion of sexual hosts in each population.** For each posterior sample, 10 simulations are run. For each population within a simulation, mean infection prevalence and mean proportion of sexual hosts is calculated by averaging across last 100 generations. Each population is assigned equal weight as the parameter that simulated the population. Colored contour lines show 95% weighted highest posterior density region. Open triangles represent observed data; proportion of sexual hosts is computed from proportion of male hosts. Open circles represent observed mean proportions averaged across years. Dotted lines around open circles represent ranges of proportion of sexual and infected hosts observed in each site.

Vergara et al. (2014b)). As transmission rate ( $\beta$ ) increases, selection for sexual hosts increases but increasing number of resistant offsprings prevents further infection from occurring and can decrease overall infection prevalence. Such a trend is consistent with previous results by Lively (2001) who noted that there is a region in which both sexual and asexual reproduction can be selected exclusively under same infection prevalence. We also find that proportion of sexual hosts decreases when infection prevalence is very high. Decrease in fitness of sexual hosts associated with increase in prevalence was predicted by Ashby and King (2015); it can also be found in an earlier work by Lively (2010b) although it was not discussed in the paper.

### Not edited:

Parameter estimates for Vergara et al. (2014b) are presented in Fig. 3. The most surprising result is that the posterior distributions of the scale parameter for cost of sex,  $c_b$ , is much wider than the prior distribution and have noticeably higher mean: (mean (95% credible interval)) 1.30 (0.80-1.97). Ashby and King (2015) defined  $c_b$  as additional costs and benefits of sex, where  $c_b = 1$  corresponds the two fold cost. Under their interpretation, our estimate corresponds to the following mean and 95% CIs for cost of sex: 1.54 (1.02-2.50). Note that some

posterior samples range  $c_b > 2$ , corresponding to faster growth rate of sexual hosts than asexual hosts. Simulated data from these posterior samples mostly consist of populations with almost 100% sexual hosts. While these estimates are not entirely impossible given that Vergara et al. (2014b) observed almost 100% sexual snails in one of their sampling sites throughout 5 years, their observation is likely to be a sampling artifact since an earlier work sampled at the same site reported much lower frequency of sexual snails (Vergara et al., 2013).

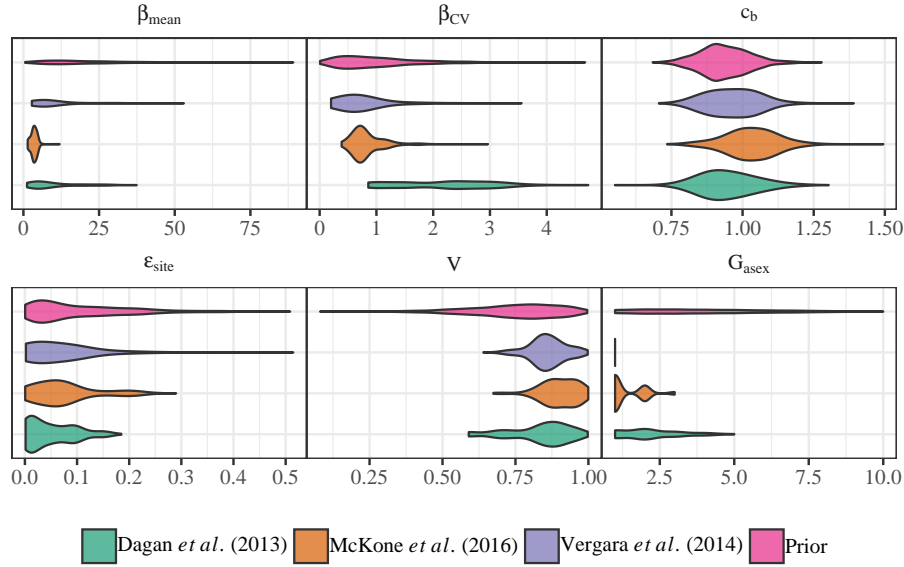


Figure 3: **Parameter estimates from Sequential Monte Carlo Approximate Bayesian Computation.** Lines represent 95% quantile and points represent mean posterior estimates. 100 posterior samples were obtained from SMC ABC. Prior distributions are specified in [TODO].

The model estimates high virulence overall: (mean (95% CI)) 87.0% (67.5%-98.7%). In contrast to Lively (2010b), who modeled 1 asexual genotype competing with 9 sexual genotypes, our model estimate shows that higher asexual to sexual genotypic ratio can be supported ( $G_{\text{asex}}$  panel in Fig. 2). 274 out of 300 posterior samples estimate  $G_{\text{asex}} = 1$ ; 25 samples estimate  $G_{\text{asex}} = 2$ ; and 1 sample estimates  $G_{\text{asex}} = 3$ . However, we find that it is still necessary for sexual hosts to have higher genetic diversity than asexual hosts.

Finally, a power analysis shows that the power for detecting a positive correlation between infection prevalence and frequency of sexual hosts in Vergara et al. (2014b) population is almost 0 (Fig. 4). We find that there is much higher (but still low) power to detect a negative correlation. Increasing number of samples per site has small effect on power once the sample size greater than 50. Increasing number of sites leads to greater increase in power but the power

appears to saturate as number of sites increases.

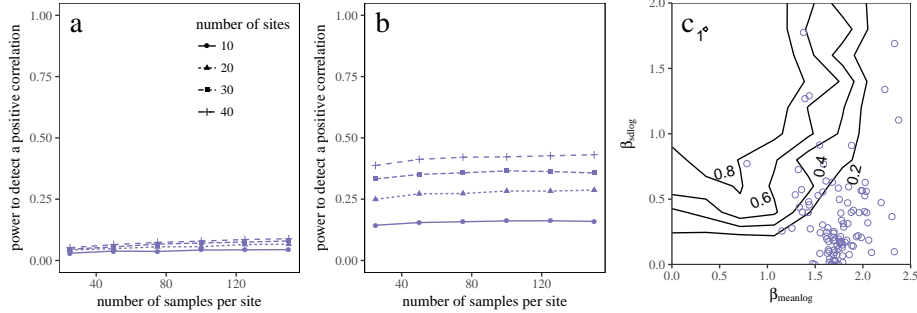


Figure 4: **Power to detect a statistically clear correlation between infection prevalence and frequency of sexual hosts.** (a) Power to detect a positive correlation. (b) Power to detect a negative correlation. (c) (Contour) power to detect a positive correlation as a function of transmission rate parameters and (points) marginal posterior distributions. Spearman’s rank correlation was used to test for correlation between infection prevalence and frequency of sexual hosts in simulated data from the posterior distributions. Contour plot is created by taking the mean posterior estimates and assuming 20 sites and 100 samples per site.

## 4 Discussion

Our results challenge ways in which the Red Queen Hypothesis for sex has been studied. Previous modeling studies have either relied on assumed parameters [CITE] or explored parameter spaces [CITE] to understand sexual reproduction maintained by host parasite coevolution. Here, we show that (1) model parameters are estimable and (2) a simple Red Queen model with population structure can reproduce key summary statistics observed across three different snail systems (Fig. 1). While our model is able to reproduce observed summary statistics, there still remain some discrepancies between model prediction and observed relationship between infection prevalence and frequency of sexual hosts across populations. These discrepancies suggest a simple host-parasite model is not sufficient to explain sexual reproduction observed in nature.

A model that does not fit well can sometimes tell us more about the system of interest than a model that fits well. For example, there is a clear mismatch between the model and the data presented by Dagan et al. (2013a) (Fig. 2). The snail populations studied by Dagan et al. (2013a) live in intrinsically different environments from two other snail populations that we consider. For example, some habitats are subject to seasonal flash floods, which can affect reproductive strategies of snails (Ben-Ami and Heller, 2007). Although we cannot estimate the strength of the effect of environment on sexual reproduction of snails relative

to that of the effect of parasites, it is likely to be a cause of the bad fit. We conclude that the Red Queen Hypothesis alone cannot explain maintenance of sexual reproduction in this population.

The only system that our model could resemble fairly well is one studied by Vergara et al. (2014b). In order for our model to produce similar summary statistics as those observed in nature, it requires the scale parameter ( $c_b$ ) for the cost of sex to have mean greater than 1 (Fig. 3). A model that assumes two fold cost of sex ( $c_b = 1$ ) would not have fitted well. However, from an inferential point of view,  $c_b > 1$  does not necessarily imply that cost of sex is less than two fold. We interpret an estimate of  $c_b > 1$  as amount of compensation required in order for the model to reproduce observed data. Additional compensation can include less than two fold cost of sex or any other mechanisms that can promote sex.

A suitable candidate for model compensation is increase in host genetic diversity. Although exact genetic architecture that determines trematode infection in snails (e.g., loci involved in parasite resistance) is not known CITE, genetic diversity of snails that have been documented is far greater than what we have assumed (King et al., 2011; Dagan et al., 2013a). In addition, increasing genetic diversity of the model may resolve overestimation of infection prevalence (Fig. 1). Previous studies have shown that moderately high genetic diversity can allow sexual hosts to escape infection more easily and therefore reduce prevalence of infection given similar amount of sexual reproduction (Lively, 2010a; King and Lively, 2012; Ashby and King, 2015). Overall, our results indicate that more modeling effort is required to understand prevalence of sexual reproduction in nature.

Our power analysis (Fig. 4) contrasts with the positive correlation predicted by Lively (1992, 2001) and findings of many empirical studies that have confirmed the prediction (Lively, 1987; Lively and Jokela, 2002; Kumpulainen et al., 2004; Vergara et al., 2013; McKone et al., 2016). The power analysis predicts almost no power for detecting a positive correlation in the population studied by Vergara et al. (2014b) and relatively higher but still low power for detecting a negative correlation. This result may appear to contradict an earlier work by Vergara et al. (2013) that reported a positive correlation between infection prevalence and male frequency in the same lake but there is a simple explanation for the difference. The key premise behind the positive correlation predicted by Lively (2001) is that there must be large variation in infection prevalence. In particular, range of prevalence must be wide enough so that the sample includes sites with almost no infected hosts (hence no sexual hosts) and those with reasonably high proportion of infected hosts to maintain sexual reproduction through parasitism (Lively, 2001). Since all four habitats studied by Vergara et al. (2014b) consists of populations with high prevalence and high frequency of sexually reproducing hosts, positive correlation vanishes. Instead, studying a system with larger variation and lower mean prevalence will yield much high power (Fig. 4(c)).

On the other hand, negative correlation between prevalence of infection and frequency of sexual hosts can be explained by cycling of host and parasite pop-

ulations. A main component of the Red Queen Hypothesis is that negative frequency dependence drives oscillation in both host and parasite population (Hamilton, 1980). When temporal variation is taken into account, association between infection prevalence and frequency of sexual hosts can change depending on what phase each of the sample population is going through in its cycle (Fig. 5). The negative correlation in the population does not contradict the positive correlation predicted by Lively (2001) because their prediction did not take temporal variation into account.

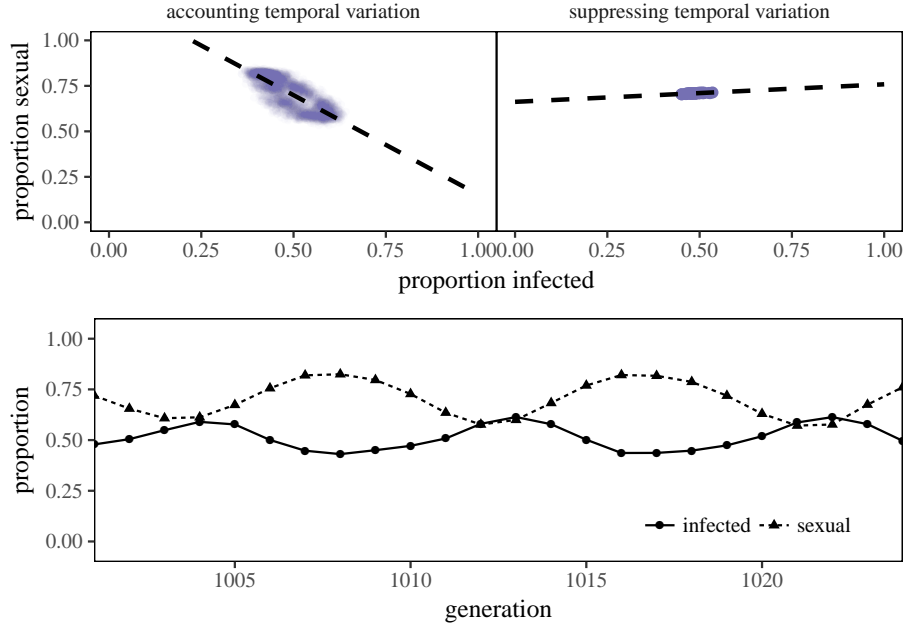


Figure 5: **Simulated data from the posterior distribution fitted to Vergara et al. (2014b).** A particular simulation is chosen from the posterior distributions to demonstrate that it is possible to predict opposite relationship when temporal variation is taken into account. (Top left) each point represents proportion of infected hosts and proportion of sexually reproducing hosts of each population at each generation. Last hundred generations are plotted. (Top right) each point represents mean proportion of infected hosts and mean proportion of sexually reproducing hosts of each population averaged over the last hundred generations. Dashed lines represent least squares fit to all point. (Bottom) A typical host parasite cycle observed in a population from this simulated data.



## References

- Agrawal, A. and C. M. Lively (2002). Infection genetics: gene-for-gene versus matching-alleles models and all points in between. *Evolutionary Ecology Research* 4(1), 91–107.
- Ashby, B. and K. C. King (2015). Diversity and the maintenance of sex by parasites. *Journal of Evolutionary Biology* 28(3), 511–520.
- Auld, S. K., S. K. Tinkler, and M. C. Tinsley (2016). Sex as a strategy against rapidly evolving parasites. *Proceedings of the Royal Society B: Biological Sciences* 283(1845).
- Bell, G. (1982). *The Masterpiece of Nature: The Evolution and Genetics of Sexuality*. University of California Press.
- Ben-Ami, F. and J. Heller (2007). Temporal patterns of geographic parthenogenesis in a freshwater snail. *Biological journal of the Linnean Society* 91(4), 711–718.
- Clarke, B. (1976). The ecological genetics of host-parasite relationships. *Genetic aspects of host-parasite relationships*. Blackwell, London, 87–103.
- Dagan, Y., K. Liljeroos, J. Jokela, and F. Ben-Ami (2013a). Clonal diversity driven by parasitism in a freshwater snail. *Journal of evolutionary biology* 26(11), 2509–2519.
- Dagan, Y., K. Liljeroos, J. Jokela, and F. Ben-Ami (2013b). Data from: Clonal diversity driven by parasitism in a freshwater snail.
- Dybdahl, M. F. and C. M. Lively (1995). Host-parasite interactions: infection of common clones in natural populations of a freshwater snail (*Potamopyrgus antipodarum*). *Proceedings of the Royal Society of London B: Biological Sciences* 260(1357), 99–103.
- Dybdahl, M. F. and C. M. Lively (1998). Host-parasite coevolution: evidence for rare advantage and time-lagged selection in a natural population. *Evolution*, 1057–1066.
- Galvani, A. P., R. M. Coleman, and N. M. Ferguson (2001). Antigenic diversity and the selective value of sex in parasites. In *Annales Zoologici Fennici*, pp. 305–314. JSTOR.
- Galvani, A. P., R. M. Coleman, and N. M. Ferguson (2003). The maintenance of sex in parasites. *Proceedings of the Royal Society of London B: Biological Sciences* 270(1510), 19–28.
- Gibson, A. K., L. F. Delph, and C. M. Lively (2017). The two-fold cost of sex: Experimental evidence from a natural system. *Evolution Letters* 1(1), 6–15.

- Gibson, A. K., J. Y. Xu, and C. M. Lively (2016). Within-population co-variation between sexual reproduction and susceptibility to local parasites. *Evolution* 70(9), 2049–2060.
- Haldane, J. B. S. (1949). Disease and evolution. *La Ricerca Scientifica Supplement* 19, 68–76.
- Hamilton, W. D. (1980). Sex versus non-sex versus parasite. *Oikos*, 282–290.
- Howard, R. S., C. M. Lively, et al. (1994). Parasitism, mutation accumulation and the maintenance of sex. *Nature* 367(6463), 554–557.
- Jaenike, J. (1978). An hypothesis to account for the maintenance of sex within populations. *Evolutionary Theory* 3, 191–194.
- Jokela, J., M. F. Dybdahl, and C. M. Lively (2009). The maintenance of sex, clonal dynamics, and host-parasite coevolution in a mixed population of sexual and asexual snails. *The American Naturalist* 174(S1), S43–S53.
- King, K. and C. M. Lively (2012). Does genetic diversity limit disease spread in natural host populations? *Heredity* 109(4), 199–203.
- King, K. C., J. Jokela, and C. M. Lively (2011). Parasites, sex, and clonal diversity in natural snail populations. *Evolution* 65(5), 1474–1481.
- Kumpulainen, T., A. Grapputo, J. Mappes, and M. Björklund (2004). Parasites and sexual reproduction in psychid moths. *Evolution* 58(7), 1511–1520.
- Lehtonen, J., M. D. Jennions, and H. Kokko (2012). The many costs of sex. *Trends in Ecology & Evolution* 27(3), 172–178.
- Lively, C. (2009). The maintenance of sex: host–parasite coevolution with density-dependent virulence. *Journal of Evolutionary Biology* 22(10), 2086–2093.
- Lively, C. M. (1987). Evidence from a new zealand snail for the maintenance of sex by parasitism. *Nature* 328(6130), 519–521.
- Lively, C. M. (1989). Adaptation by a parasitic trematode to local populations of its snail host. *Evolution* 43(8), 1663–1671.
- Lively, C. M. (1992). Parthenogenesis in a freshwater snail: reproductive assurance versus parasitic release. *Evolution* 46(4), 907–913.
- Lively, C. M. (2001). Trematode infection and the distribution and dynamics of parthenogenetic snail populations. *Parasitology* 123(07), 19–26.
- Lively, C. M. (2010a). The effect of host genetic diversity on disease spread. *The American Naturalist* 175(6), E149–E152.
- Lively, C. M. (2010b). An epidemiological model of host–parasite coevolution and sex. *Journal of evolutionary biology* 23(7), 1490–1497.

- Lively, C. M. and J. Jokela (2002). Temporal and spatial distributions of parasites and sex in a freshwater snail. *Evolutionary Ecology Research* 4(2), 219–226.
- MacPherson, A. and S. P. Otto (2017). Joint coevolutionary-epidemiological models dampen red queen cycles and alter conditions for epidemics. *Theoretical population biology*.
- May, R. M. and R. M. Anderson (1983). Epidemiology and genetics in the coevolution of parasites and hosts. *Proceedings of the Royal Society of London B: Biological Sciences* 219(1216), 281–313.
- McKone, M. J., A. K. Gibson, D. Cook, L. A. Freymiller, D. Mishkind, A. Quinlan, J. M. York, C. M. Lively, and M. Neiman (2016). Fine-scale association between parasites and sex in *Potamopyrgus antipodarum* within a New Zealand lake. *New Zealand Journal of Ecology* 40(3), 1.
- Morris, C. N. et al. (1983). Natural exponential families with quadratic variance functions: statistical theory. *The Annals of Statistics* 11(2), 515–529.
- Otto, S. P. (2009). The evolutionary enigma of sex. *The American naturalist* 174(S1), S1–S14.
- Otto, S. P. and Y. Michalakis (1998). The evolution of recombination in changing environments. *Trends in Ecology & Evolution* 13(4), 145–151.
- Slowinski, S. P., L. T. Morran, R. C. Parrish, E. R. Cui, A. Bhattacharya, C. M. Lively, and P. C. Phillips (2016). Coevolutionary interactions with parasites constrain the spread of self-fertilization into outcrossing host populations. *Evolution* 70(11), 2632–2639.
- Smith, J. M. (1978). *The Evolution of Sex*, Volume 54. Cambridge Univ Press.
- Smith, J. M. and M. Slatkin (1973). The stability of predator-prey systems. *Ecology* 54(2), 384–391.
- Tobler, M. and I. Schlupp (2008). Expanding the horizon: the red queen and potential alternatives. *Canadian Journal of Zoology* 86(8), 765–773.
- Toni, T., D. Welch, N. Strelkowa, A. Ipsen, and M. P. Stumpf (2009). Approximate bayesian computation scheme for parameter inference and model selection in dynamical systems. *Journal of the Royal Society Interface* 6(31), 187–202.
- Turner, B. M. and T. Van Zandt (2012). A tutorial on approximate bayesian computation. *Journal of Mathematical Psychology* 56(2), 69–85.
- Vergara, D., J. Jokela, and C. Lively (2014a). Data from: Infection dynamics in coexisting sexual and asexual host populations: support for the red queen hypothesis.

- Vergara, D., J. Jokela, and C. M. Lively (2014b). Infection dynamics in co-existing sexual and asexual host populations: support for the Red Queen hypothesis. *The American naturalist* 184(S1), S22–S30.
- Vergara, D., C. M. Lively, K. C. King, and J. Jokela (2013). The geographic mosaic of sex and infection in lake populations of a new zealand snail at multiple spatial scales. *The American Naturalist* 182(4), 484–493.
- Vrijenhoek, R. C. (1998). Animal clones and diversity. *Bioscience* 48(8), 617–628.