

INFORMING POLICY VIA DYNAMIC MODELS: ELIMINATING CHOLERA IN HAITI

BY JESSE WHEELER^{1,a} , ANNAELAINE ROSENGART^{1,b} ZHUOXUN JIANG^{1,c},
KEVIN HAO EN TAN^{1,d}, NOAH TREUTLE^{1,e} AND EDWARD IONIDES^{1,f}

¹STATISTICS DEPARTMENT, UNIVERSITY OF MICHIGAN, ^aJESWHEEL@UMICH.EDU; ^bAELR@UMICH.EDU;
^cZHUOXUNJ@UMICH.EDU; ^dKEVTAN@UMICH.EDU; ^eNTREUTLE@UMICH.EDU; ^fIONIDES@UMICH.EDU

Working draft, October 17, 2022

Public health decisions must be made about when and how to implement interventions to control an infectious disease epidemic. These decisions should be informed by data on the epidemic as well as current understanding about the transmission dynamics. Such decisions can be posed as statistical questions about scientifically motivated dynamic models. Thus, we encounter the methodological task of building credible, data-informed decisions based on stochastic, partially observed, nonlinear dynamic models. This necessitates addressing the tradeoff between biological fidelity and model simplicity, and the reality of misspecification for models at all levels of complexity. As a case study, we consider a cholera epidemic in Haiti. The 2010 introduction of cholera to Haiti led to an extensive outbreak and sustained transmission until it was eliminated in 2019. We study three models developed by expert teams to advise on vaccination policies. We assess methods used for fitting and evaluating these models, leading to recommendations for future studies. Diagnosis of model misspecification and development of alternative models can lead to improved statistical fit, but caution is nevertheless required in drawing policy conclusions based on causal interpretations of the models.

1. Introduction. Regulation of biological populations is a fundamental topic in epidemiology, ecology, fisheries and agriculture. Population dynamics may be nonlinear and stochastic, with the resulting complexities compounded by incomplete understanding of the underlying biological mechanisms and by partial observability of the system variables. Quantitative models for these dynamic systems offer potential for designing effective control measures. Developing and testing such models, and assessing their fitness for guiding policy, is a challenging statistical task. Questions of interest include: What indications should we look for in the data to assess whether the model-based inferences are trustworthy? What diagnostic tests and model variations can and should be considered in the course of the data analysis? What are the possible trade-offs of increasing model complexity, such as the inclusion of interactions across spatial units?

This case study investigates the use of dynamic models and spatiotemporal data to inform a policy decision in the context of the cholera outbreak in Haiti, which started in 2010. We build on a multi-group modeling exercise by Lee et al. (2020a) in which four expert modeling teams developed models to the same dataset with the goal of comparing conclusions on the feasibility of eliminating cholera by a vaccination campaign. Model 1 is stochastic and describes cholera at the national level; Model 2 is

Keywords and phrases: Partially observed Markov process, Hidden Markov model, infectious disease, cholera, sequential Monte Carlo.

deterministic with spatial structure, and includes transmission via contaminated water; Model 3 is stochastic with spatial structure, and accounts for measured rainfall. Model 4 has an agent-based construction, featuring considerable mechanistic detail but limited ability to calibrate these details to data. The strengths and weaknesses of the agent-based modeling approach (Tracy, Cerdá and Keyes, 2018) are outside the scope of this article, and we focus on Models 1–3.

The four independent teams were given the task of estimating the potential effect of prospective oral cholera vaccine (OCV) programs. While OCV is accepted as a safe and effective tool for controlling the spread of cholera, the global stockpile of OCV doses remains limited (Pezzoli, 2020). Advances in OCV technology and vaccine availability, however, raised the possibility of planning a national vaccination program (Lee et al., 2020a). In the study, certain data were shared between the groups, including demography and vaccination history; vaccine efficacy was also fixed at a shared value between groups. Beyond this, the groups made autonomous decisions on what to include and exclude from their models. The groups largely adhered to existing guidelines on creating models to inform policy (Behrend et al., 2020; Saltelli et al., 2020) and, despite their autonomy, obtained a consensus that an extensive nationwide vaccination campaign would be necessary to eliminate cholera from Haiti. This conclusion is inconsistent with the fact that there have been no confirmed cases since February, 2019 (Trevisin et al., 2022) without the implementation of a national vaccination program.

The discrepancy between the model based conclusions of Lee et al. (2020a) and the subsequent elimination of cholera has been debated (Francois, 2020; Rebaudet, Gaudart and Piarroux, 2020; Henrys et al., 2020; Lee et al., 2020b). Suggested origins of this discrepancy include the use of unrealistic models (Rebaudet, Gaudart and Piarroux, 2020) and unrealistic criteria for cholera elimination (Henrys et al., 2020). We find a more nuanced conclusion: attention to methodological details in model fitting, diagnosis and forecasting can improve each of the proposed model’s ability to quantitatively describe observed data. These improvements result in forecasts that are more consistent with the observed outcome, without requiring major changes to the model structures. Based on this retrospective analysis, we offer suggestions on fitting mechanistic models to dynamic systems for future studies.

While this work is focused on the models proposed by Lee et al. (2020a), our suggestions have broader relevance. To investigate the extent to which Lee et al. (2020a) is typical of the substantial body of work on Haiti cholera dynamics, we performed a literature review by searching PubMed with keywords: Haiti, cholera, model. The search resulted in 66 papers, of which 32 used dynamic models to describe the cholera epidemic in Haiti. The models make various choices on the dichotomies considered by Lee et al. (2020a); deterministic versus stochastic; compartment model versus agent-based; aggregated versus spatially explicit. This is no accident, since Lee et al. (2020a) purposefully designed their study to cover the range of current modeling practice.

We proceed by introducing Models 1–3 in Sec. 2; in Sec. 3, we present a systematic approach to examining and refining these models, and then use improved model fits to project cholera incidence in Haiti under various vaccination scenarios. This is followed by a consideration of the robustness of model based policy recommendations in Sec. 4, and a discussion in Sec. 5.

2. Mechanistic models for cholera in Haiti. Models that focus on learning relationships between variables in a dataset are called *associative*, whereas models that incorporate a known scientific property of the system are called *causal* or *mechanistic*. The danger in using forecasting techniques which rely on associative models to predict



FIG 1. Weekly reported cholera cases in Haiti from October 2010 to January 2019.

the consequence of interventions is called the Lucas critique in an econometric context. Lucas et al. (1976) pointed out that it is naive to predict the effects of an intervention on a given system based entirely on historical associations. To successfully predict the effect of an intervention, a model should therefore both provide a quantitative explanation of existing data and should have a causal interpretation: a manipulation of the system should correspond quantitatively with the corresponding change to the model. This motivates the development of mechanistic statistical models, which provides a statistical fit to the available data while also supporting a causal interpretation.

The four mechanistic models of Lee et al. (2020a) were deliberately developed with limited coordination. This allows us to treat the models as fairly independently developed expert approaches to understanding cholera transmission. However, it led to differences in notation, and in subsets of the data chosen for analysis, that hinder direct comparison. Here, we have put all three models into a common notational framework. Translations back to the original notation of Lee et al. (2020a) are given in Table S-1.

Each model describes the cholera dynamics as a partially observed Markov process (POMP) with a latent state vector $\mathbf{X}^{(m)}(t)$ for each continuous time point t and model $m \in \{1, 2, 3\}$. N observations on the system are collected at time points t_1, \dots, t_N , written as $t_{1:N}$. The observation at time t_n is modeled by the random vector $\mathbf{Y}_n^{(m)}$. While the latent process exists between observation times, the value of the latent state at observations times is of particular interest. We therefore write $\mathbf{X}_n^{(m)} = \mathbf{X}^{(m)}(t_n)$ to denote the value of the latent process at the n th observation time, and $\mathbf{X}_{1:N}^{(m)}$ is the collection of latent state values for all observed time points. The observable random variables $\mathbf{Y}_{1:N}^{(m)}$ are assumed to be conditionally independent given $\mathbf{X}_{0:N}^{(m)}$. Together, with the density for the initial value of the latent state $\mathbf{X}_0^{(m)} = \mathbf{X}^{(m)}(t_0)$, each model assumes a joint density $f_{\mathbf{X}_{0:N}^{(m)}, \mathbf{Y}_{1:N}^{(m)}}(\mathbf{x}_{0:N}^{(m)}, \mathbf{y}_{1:N}^{(m)}; \theta)$, where θ is a parameter vector that indexes the model. The observed data $\mathbf{y}_{1:N}^*$, along with the unobserved true value of the latent state, are modeled as a realization of this joint distribution.

For each model $m \in \{1, 2, 3\}$, the latent state vector $\mathbf{X}^{(m)}(t)$ consists of individuals labeled as susceptible (S), infected (I), asymptotically infected (A), vaccinated (V),

and recovered (R), with various sub-divisions sometimes considered. Models 2 and 3 have metapopulation structure, meaning that each individual is a member of a spatial unit, denoted by a subscript $u \in 1:U$. Here, the spatial units are the $U = 10$ Haitian administrative départements (henceforth anglicized as departments).

In the following subsections, descriptions of Models 1–3 are provided. While the model description is scientifically critical, as well as necessary for transparency and reproducibility, the model details are not essential to our methodological discussions of how to diagnose and address model misspecification with the purpose of informing policy. A first-time reader may choose to skim through the rest of this section, and return later.

2.1. *Model 1.* $\mathbf{X}^{(1)}(t) = (S_z(t), E_z(t), I_z(t), A_z(t), R_z(t), z \in 0:Z)$ describes susceptible, latent (exposed), infected (and symptomatic), asymptomatic, and recovered individuals in vaccine cohort z . Here, $z = 0$ corresponds to unvaccinated individuals, and $z \in 1:Z$ describes hypothetical vaccination programs. The force of infection is

$$(1) \quad \lambda(t) = \left(\sum_{z=0}^Z I_z(t) + \epsilon \sum_{z=0}^Z A_z(t) \right)^\nu \frac{d\Gamma(t)}{dt} \beta(t) / N,$$

where $\beta(t)$ is a periodic cubic spline representation of seasonality, given in terms of a B-spline basis $\{s_j(t), j \in 1:6\}$ and parameters $\beta_{1:6}$ as

$$(2) \quad \log \beta(t) = \sum_{j=1}^6 \beta_j s_j(t).$$

The process noise $d\Gamma(t)/dt$ is multiplicative Gamma-distributed white noise, with infinitesimal variance parameter σ_{proc}^2 . Lee et al. (2020a) included process noise in Model 3 but not in Model 1, i.e., they fixed $\sigma_{\text{proc}}^2 = 0$. Gamma white noise in the transmission rate gives rise to an over-dispersed latent Markov process (Bretó and Ionides, 2011) which has been found to improve the statistical fit of disease transmission models (Stocks, Britton and Höhle, 2020; He, Ionides and King, 2010).

Per-capita transition rates are given in Equations 3-10:

$$\begin{aligned} (3) \quad & \mu_{S_z E_z} = \lambda(t), \\ (4) \quad & \mu_{E_z I_z} = \mu_{EI} (1 - f_z(t)), \\ (5) \quad & \mu_{E_z A_z} = \mu_{EI} f_z(t), \\ (6) \quad & \mu_{I_z R_z} = \mu_{A_z R_z} = \mu_{IR}, \\ (7) \quad & \mu_{R_z S_z} = \mu_{RS}, \\ (8) \quad & \mu_{S_0 S_z} = \mu_{E_0 E_z} = \mu_{I_0 I_z} = \mu_{A_0 A_z} = \mu_{R_0 R_z} = \eta_z(t), \\ (9) \quad & \mu_{S_z D} = \mu_{E_z D} = \mu_{I_z D} = \mu_{A_z D} = \mu_{R_z D} = \delta, \\ (10) \quad & \mu_{DS_0} = \mu_S, \end{aligned}$$

where $z \in 0:Z$. Here, μ_{AB} is a transition rate from compartment A to B . We have an additional demographic source and sink compartment D modeling entry into the study population due to birth or immigration, and exit from the study population due to death or immigration. Thus, μ_{AD} is a rate of exiting the study population from compartment A and μ_{DB} is a rate of entering the study population into compartment B .

In Model 1, the advantage afforded to vaccinated individuals is an increased probability that an infection is asymptomatic. Conditional on infection status, vaccinated

individuals are also less infectious than their non-vaccinated counterparts by a rate of $\epsilon = 0.05$ in Eq. (1). In (5) and (4) the asymptomatic ratio for non-vaccinated individuals is set $f_0(t) = 0$, so that the asymptomatic route is reserved for vaccinated individuals. For $z \in 1:Z$, the vaccination cohort z is assigned a time τ_z , and we take $f_z(t) = c\theta^*(t - \tau_z)$ where $\theta^*(t)$ is efficacy at time t since vaccination for adults, taken from Lee et al. (2020a), Table S4, and $c = (1 - (1 - 0.4688) \times 0.11)$ is a correction to allow for reduced efficacy in the 11% of the population aged under 5 years. Single and double vaccine doses were modeled by changing the waning of protection; protection was assumed to be equal between single and double dose until 52 weeks after vaccination, at which point the single dose becomes ineffective. The latent state vector $\mathbf{X}^{(1)}(t)$ is initialized by setting each compartment count for each vaccination scenario $z \neq 0$ as zero, and introducing initial-value parameters $I_{0,0}$ and $E_{0,0}$ such that $R_0(0) = 0$, $I_0(0) = \text{Pop} \times I_{0,0}$, $E_0(0) = \text{Pop} \times E_{0,0}$ and $S_0(0) = \text{Pop} \times (1 - I_{0,0} - E_{0,0})$, where Pop is the total population of Haiti. Reported cholera cases at time point n ($\mathbf{Y}_n^{(1)}$) are assumed to come from a negative binomial measurement model, where only a fraction (ρ) of new weekly cases are reported. See Sec. S4 of the supplement material for more details.

2.2. Model 2. Susceptible individuals are in compartments $S_{uz}(t)$, where $u \in 1:U$ corresponds to the $U = 10$ departments, and $z \in 0:4$ describes vaccination status:

- $z = 0$: Unvaccinated or waned vaccination protection.
- $z = 1$: One dose at age under five years.
- $z = 2$: Two doses at age under five years.
- $z = 3$: One dose at age over five years.
- $z = 4$: Two doses at age over five years.

Individuals can progress to a latent infection E_{uz} followed by symptomatic infection I_{uz} with recovery to R_{uz} or asymptomatic infection A_{uz} with recovery to R_{uz}^A . The force of infection depends on both direct transmission and an aquatic reservoir, $W_u(t)$, and is given by

$$(11) \quad \lambda_u(t) = 0.5(1 + a \cos(2\pi t + \phi)) \frac{\beta_W W_u(t)}{W_{\text{sat}} + W_u(t)} + \beta \left\{ \sum_{z=0}^4 I_{uz}(t) + \epsilon \sum_{z=0}^4 A_{uz}(t) \right\}.$$

The latent state is therefore described by the vector $\mathbf{X}^{(2)}(t) = (S_{uz}(t), E_{uz}(t), I_{uz}(t), A_{uz}(t), R_{uz}(t), R_{uz}^A(t), W_u, u \in 1:U, z \in 0:4)$. The cosine term in Eq. (11) accounts for annual seasonality, with a phase parameter ϕ . The Lee et al. (2020a) implementation of Model 2 fixes $\phi = 0$.

Individuals move from department u to v at rate T_{uv} , and aquatic cholera moves at rate T_{uv}^W . The nonzero transition rates are

$$(12) \quad \mu_{S_{uz}E_{uz}} = \theta_z \lambda,$$

$$(13) \quad \mu_{E_{uz}I_{uz}} = f\mu_{EI}, \quad \mu_{E_{uz}A_{uz}} = (1 - f)\mu_{EI},$$

$$(14) \quad \mu_{I_{uz}R_{uz}} = \mu_{A_{uz}R_{uz}^A} = \mu_{IR},$$

$$(15) \quad \mu_{R_{uz}S_{uz}} = \mu_{R_{uz}^A S_{uz}} = \mu_{RS},$$

$$(16) \quad \mu_{S_{uz}S_{vz}} = \mu_{E_{uz}E_{vz}} = \mu_{I_{uz}I_{vz}} = \mu_{A_{uz}A_{vz}} = \mu_{R_{uz}R_{vz}} = \mu_{R_{uz}^A R_{vz}^A} = T_{uv},$$

$$(17) \quad \mu_{S_{u1}S_{u0}} = \mu_{S_{u3}S_{u0}} = \omega_1,$$

$$(18) \quad \mu_{S_{u2}S_{u0}} = \mu_{S_{u4}S_{u0}} = \omega_2,$$

$$(19) \quad \mu_{DW_u} = \mu_W \left\{ \sum_{z=0}^4 I_{uz}(t) + \epsilon_W \sum_{z=0}^4 A_{uz}(t) \right\},$$

$$(20) \quad \mu_{W_u D} = \delta_W,$$

$$(21) \quad \mu_{W_u W_v} = w_r T_{uv}^W.$$

In (16) the spatial coupling is specified by a gravity model,

$$(22) \quad T_{uv} = v_{\text{rate}} \times \frac{\text{Pop}_u \text{Pop}_v}{D_{uv}^2},$$

where Pop_u is the mean population for department u , D_{uv} is a distance measure estimating average road distance between randomly chosen members of each population, and $v_{\text{rate}} = 10^{-12}$ was treated as a fixed constant. In (21), T_{uv}^W is a measure of river flow between departments. The unit of $W_u(t)$ is cells per ml, with dose response modeled via a saturation constant of W_{sat} in (11). The starting value for each element of the latent state vector $\mathbf{X}^{(2)}(0)$ are set to zero except for $I_{u0}(0) = y_u^*(0)/\rho$ and $R_{u0}(0) = \text{Pop}_u - I_{u0}(0)$, where $y_u^*(0)$ is the reported number of cholera cases in department u at time $t = 0$. Reported cases are assumed to come from a log-normal distribution, with the log-scale mean equal to the reporting rate ρ times the number of newly infected individuals. See Sec. 3.1.2 and Sec. S4.2 of the supplement for more details.

2.3. Model 3. The latent state is described as $\mathbf{X}^{(3)}(t) = (S_{uz}(t), I_{uz}(t), A_{uz}(t), R_{uzk}(t), W_u(t), u \in 0:U, z \in 0:4, k \in 1:3)$. Here, $z = 0$ corresponds to unvaccinated, $z = 2j - 1$ corresponds to a single dose on the j th vaccination campaign in unit u and $z = 2j$ corresponds to receiving two doses on the j th vaccination campaign. $k \in 1:3$ models non-exponential duration in the recovered class before waning of immunity. The force of infection is

$$(23) \quad \lambda_u(t) = \beta_{W_u} \frac{W_u(t)}{1 + W_u(t)} + \beta_u \sum_{v \neq u} (I_{v0}(t) + \epsilon A_{v0}(t)),$$

and Per-capita transition rates are given in Equations 24–32.

$$(24) \quad \mu_{S_{uz} I_{uz}} = f \lambda_u (1 - \eta_{uz}(t)) d\Gamma/dt,$$

$$(25) \quad \mu_{S_{uz} A_{uz}} = (1 - f) \lambda_u (1 - \eta_{uz}(t)) d\Gamma/dt,$$

$$(26) \quad \mu_{I_{uz} R_{uz1}} = \mu_{A_{uz} R_{uz1}} = \mu_{IR},$$

$$(27) \quad \mu_{I_{uz} S_{u0}} = \delta + \delta_C,$$

$$(28) \quad \mu_{A_{uz} S_{u0}} = \delta,$$

$$(29) \quad \mu_{R_{uz1} R_{uz2}} = \mu_{R_{uz2} R_{uz3}} = 3\mu_{RS},$$

$$(30) \quad \mu_{R_{uzk} S_{u0}} = \delta + 3\mu_{RS} \mathbf{1}_{\{k=3\}},$$

$$(31) \quad \mu_{DW_u} = [1 + a(J_u(t))^r] \text{Den}_u \mu_W [I_{u0}(t) + \epsilon_W A_{u0}(t)],$$

$$(32) \quad \mu_{W_u D} = \delta_W.$$

As with Model 1, $d\Gamma_u(t)/dt$ is multiplicative Gamma-distributed white noise in (24) and (25). In (31), $J_u(t)$ is a dimensionless measurement of precipitation that has been standardized by dividing the observed rainfall at time t by the maximum recorded rainfall in department u during the epidemic, and Den_u is the population density. Demographic stochasticity is accounted for by modeling non-cholera related death rate δ in each compartment, along with an additional death rate δ_C in (27) to account for cholera induced deaths among infected individuals. All deaths are balanced by births into the susceptible compartment in (28) and (30), thereby maintaining constant population in each department.

Latent states are initialized using an approximation of the instantaneous number of infected, asymptomatic, and recovered individuals at time $t = 0$ by using the first week of cholera incidence data. Specifically, we set $I_{u0}(0) = \frac{y_{u0}^*}{\rho(\delta + \delta_C + \mu_{IR})}$, $A_{u0}(0) = \frac{1-f}{f} I_{u0}(0)$, and $R_{u0k} = y_{u0}^* - I_{u0}(0) - A_{u0}(0)$, and all other compartments that represent population counts are set to zero at time $t = 0$. We replace $J_u(0)$ in (31) with initial value parameters ξ_u to allow for some flexibility in initializing the bacterial compartment $W_u(0)$ while still maintaining the primary dynamics of the model. Reported cholera cases are assumed to come from a negative binomial measurement model with mean equal to a fraction (ρ) of individuals in each unit who develop symptoms seek healthcare (see Sec. S4 of the supplement material for more details).

3. Statistical Analysis. We consider model fitting (Sec. 3.1) followed by diagnostic investigations (Sec. 3.2), forecasting (Sec. 3.3) and external scientific corroboration of model fit (Sec. 3.4).

3.1. Model Fitting. Each of the three models considered in this study describes cholera dynamics as a partially observed Markov process (POMP), with the understanding that the deterministic Model 2 is a degenerate case of a stochastic model. Each model is indexed by a parameter vector, θ , and different values of θ can result in qualitative differences in the predicted behavior of the system. Therefore, the choice of θ used to make inference about the system can greatly affect model based conclusions. Elements of θ can be fixed at a constant value based on scientific understanding of the system, but parameters can also be calibrated to data by maximizing a measure of congruency between the observed data and the assumed mechanistic structure. Calibrating model parameters does not guarantee that the resulting model successfully approximates real-world mechanisms, since model assumptions may be incorrect, and do not change as the model is calibrated to data. However, the congruency between the model and observed data serves as a proxy for the congruency between the model and the true underlying dynamic system, since the observed data are a result of the system. In the following subsections we describe our approach to fitting the three proposed mechanistic models by calibrating θ to observed cholera incidence data.

3.1.1. Calibrating Model 1 Parameters. Model 1 describes cholera dynamics at the nationally aggregated scale so that $\mathbf{Y}_n^{(1)} \in \mathbb{R}$ for each observation time $n \in 1 : N$, and the latent state vector $\mathbf{X}^{(1)}(t)$ is comprised of national level population counts. Several algorithms exist, both frequentist and Bayesian, that can be used to obtain estimates of the parameters for the class of models. In order to retain the ability to propose models that are scientifically meaningful rather than only those that are simply statistically convenient, we restrict ourselves to parameter estimation techniques that have the plug-and-play property, which is that the fitting procedure only requires the ability to simulate the latent process instead of evaluating transition densities (Bretó et al., 2009; He, Ionides and King, 2010). Plug-and-play algorithms include Bayesian approaches like ABC and PMCMC (Toni et al., 2009; Andrieu, Doucet and Holenstein, 2010), but here we use frequentist methods to maximize model likelihoods. To our knowledge, the only plug-and-play frequentist methods that can maximize the likelihood for POMP models of this complexity are iterated filtering algorithms, which modify the well-known particle filter (Arulampalam et al., 2002) by performing a random walk for each parameter and particle. These perturbations are carried out iteratively over multiple filtering operations, using the collection of parameters from the previous filtering pass as the

parameter initialization for the next iteration, and decreasing the random walk variance at each step.

The ability to maximize the likelihood allows for likelihood-based inference, such as performing statistical tests for potential model improvements. We demonstrate this capability by proposing a linear trend ζ in transmission in Eq. (2):

$$(33) \quad \log \beta(t) = \sum_{j=1}^6 \beta_s s_j(t) + \zeta \bar{t},$$

where $\bar{t} = \frac{t - (t_N + t_0)/2}{t_N - (t_N + t_0)/2}$, so that $\bar{t} \in [-1, 1]$. The proposal of a linear trend in transmission is a result of observing an apparent decrease in reported cholera infections from 2012–2019 in Fig. 1. While several factors may contribute to this decrease, one explanation is that case-area targeted interventions (CATIs), which included education sessions, increased monitoring, household decontamination, soap distribution, and water chlorination in infected areas (Rebaudet et al., 2019), may have greatly reduced cholera transmission (Rebaudet et al., 2021).

We perform a statistical test to determine whether or not the data indicate the presence of a linear trend in transmissibility. To do this, we perform a profile-likelihood search on the parameter ζ and obtain a confidence interval via a Monte Carlo Adjusted Profile (MCAP) (Ionides et al., 2017). Lee et al. (2020a) implemented Model 1 by fitting two distinct phases: an epidemic phase from October 2010 through March 2015, and an endemic phase from March 2015 onward. We similarly allow the re-estimation of process and measurement overdispersion parameters (σ_{proc}^2 and ψ), and require that the latent Markov process $X(t)$ carry over from one phase into the next. The resulting confidence interval for ζ is $(-0.085, -0.005)$, with the full results displayed in Fig. 2. These results are suggestive that the inclusion of a trend in transmission rate improves the quantitative ability of Model 1 to describe the observed data. The reported results for Model 1 in the remainder of this article were obtained with the inclusion of the parameter ζ .

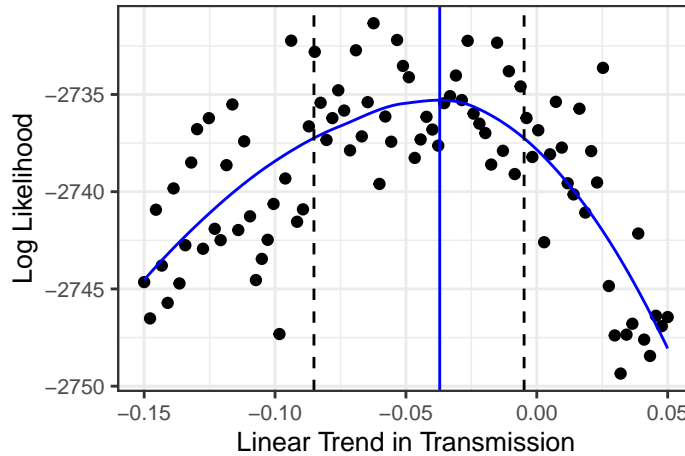


FIG 2. Monte Carlo adjusted profile of ζ . The blue curve is the profile, the blue line indicates the MLE, and the dashed lines indicate the confidence interval.

We implemented Model 1 using the `pomp` package King et al. (2009), relying heavily on the `pomp` source code provided by Lee et al. (2020a). Both analyses used the `mif2`

implementation of the IF2 algorithm to estimate θ by maximum likelihood. One change we made in the statistical analysis that led to larger model likelihoods was increasing the computational effort in the numerical maximization. While IF2 is a powerful method that can be used to fit parameters for a large class of models, the theoretic ability to maximize the likelihood depends on asymptotics in both the number of particles and the number of filtering iterations and therefore requires a large number of computations. The large difference in likelihoods that were obtained by these two approaches (see Table 1) highlights the importance of carefully determining the necessary computational effort needed to maximize model likelihoods and acting accordingly.

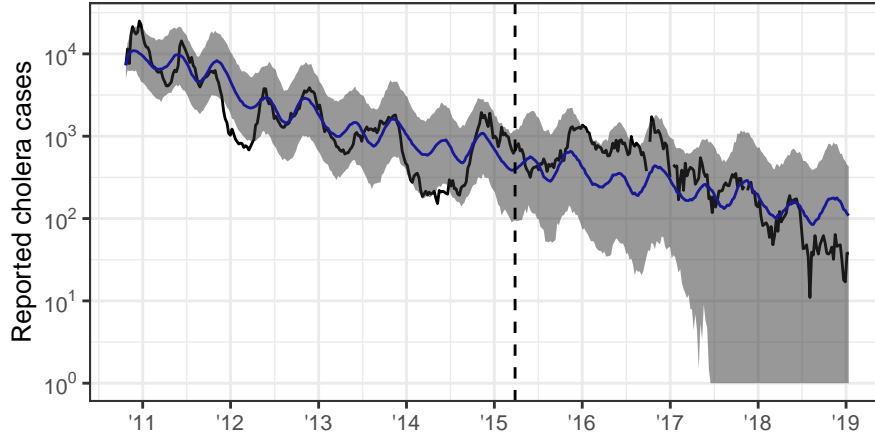


FIG 3. Simulations from Model 1 compared to reported cholera cases. The black curve is observed data, the blue curve is median of 500 simulations from the fitted model, and the vertical dashed line represents break-point when parameters are refit.

3.1.2. Calibrating Model 2 Parameters. As previously stated, Model 2 is a deterministic compartmental model defined by a set of coupled differential equations. The use of deterministic compartment models have a long history in the field of infectious disease epidemiology (Kermack and McKendrick, 1927; Brauer, 2017; Giordano et al., 2020), and can be justified by asymptotic considerations in a large-population limit (Dadlani et al., 2020; Ndii and Supriatna, 2017). Because the process model of Model 2 is deterministic, the parameter estimation problem for Model 2 reduces to a least squares calculation when combined with a Gaussian measurement model (see supplement material for details). Lee et al. (2020a) fit two versions of model 2 based on a presupposed change in cholera transmission from an epidemic phase to endemic phase that occurred in March, 2014. The inclusion of a change-point in model states and parameters increases the flexibility of the model and hence the ability to fit the observed data. The increase in model flexibility, however, results in hidden states that are inconsistent between model phases. The inclusion of a model break-point by Lee et al. (2020a) is perhaps due to a challenging feature of fitting a deterministic model via least squares: discrepancies between model trajectories and observed case counts in highly infectious periods of a disease outbreak will result in greater penalty than the discrepancies between model trajectories and observed case counts in times of relatively low infectiousness, resulting in a bias towards accurately describing periods of high infectiousness. This issue is particularly troublesome for modeling cholera dynamics in

Haiti: the inability to accurately fit times of low infectiousness may result in poor model forecasts, as few cases of cholera were observed the last few years of the epidemic.

To combat this issue, we fit the model to log-transformed case counts, since the log scale stabilizes the variation during periods of high and low incidence. Other solutions include changing the measurement model to include overdispersion, as was done in Models 1 and 3. This permits the consideration of demographic stochasticity, which is dominant for small infected populations, together with log scale stochasticity (also called multiplicative, or environmental, or extra-demographic) which is dominant at high population counts. Here we chose to fit the model to transformed case counts rather than adding overdispersion to the measurement model with the goal of making the fewest changes to the model proposed by Lee et al. (2020a) as possible.

We implemented this model using the `spatPomp` R package (Asfaw, Ionides and King, 2021). The model was then fit using the subplex algorithm, implemented in the `subplex` package (King and Rowan, 2020). A comparison of the trajectory of the fitted model to the data is given in Fig 6.

3.1.3. Calibrating Model 3 Parameters. Both the latent and observable processes in Model 3 can be factored into department specific processes which interact with each other. The decision to address metapopulation dynamics via a spatially explicit model, rather than to aggregate over space, is double-edged. Evidence for the former approach has been provided in previous studies (King et al., 2015), including the specific case of heterogeneity between Haitian departments in cholera transmission (Collins and Govinder, 2014). However, a legitimate preference for simplicity (Saltelli et al., 2020; Green and Armstrong, 2015) can support a decision to consider nationally aggregated models.

Fitting scientifically flexible metapopulation models is a challenging statistical problem. In particular, parameter estimation techniques based on the particle filter become computationally intractable as the number of spatial units increase. This is a result of the approximation error of particle filters growing exponentially in the dimension of the model (Rebeschini and van Handel, 2015; Park and Ionides, 2020).

Model 3 parameters that are calibrated to data are primarily shared between each department, meaning that the estimated parameter value for department u is the same value as the parameter value in department v , with $u, v \in 1:U$. The exception to this being the parameters β_{W_u} , and β_u , which take department-specific values for each department u . To avoid the parameter estimation issue in high-dimensional models, Lee et al. (2020a) simplified the problem by using the IF2 algorithm to estimate parameters separately for independent department-level models; the shared parameters were calibrated using the cholera incidence data from Artibonite, and the department-specific parameters (β_{W_u} and β_u) were fit using the data from their respective department. Reducing a spatially coupled model to individual units in this fashion requires special treatment of any interactive mechanisms between spatial units, such as found in Eq. (23). In particular, when considering a model for department u , the values $I_\nu(t)$ and $A_\nu(t)$, are unknown for $u \neq v$. To address this issue, Lee et al. (2020a) conditioned each department level model on reported cholera cases in other departments and approximated the values of $I_\nu(t)$ and $A_\nu(t)$ using the respective reported cases.

The simplified, spatially-decoupled version of Model 3 relies on the observed cholera cases for each department to describe cholera dynamics, making it impossible to obtain forecasts for future time points from the fitted model. In order to obtain model forecasts, the parameters fit for the conditionally independent department level models were used in the fully coupled version of Model 3. This approach of model calibration and

forecasting avoids the issue of particle depletion by simplifying the parameter estimation problem, but this advantage raises new issues. One concern is that cholera dynamics in department u are highly related to the dynamics in the remaining departments; calibrating model parameters using observed cases in other departments as a model covariate may therefore lead to an over-dependence on observed cholera cases. Another concern is that the two versions of the model are not the same, resulting in sub-optimal parameter estimates for the spatially coupled model. This is because the likelihood of the coupled model at its maximum likelihood estimate (MLE) is greater than or equal to the likelihood of the coupled model at the MLE of the decoupled model. These two concerns may explain the unrealistic forecasts and low likelihood of the model (see Table 1).

For our analysis, we calibrate the parameters of the spatially coupled version of Model 3 using the iterated block particle filter (IBPF) algorithm of Ionides, Ning and Wheeler (2022). This algorithm extends the work of Ning and Ionides (2021), who provided theoretic justification for the version of the algorithm that only estimates unit-specific parameters. While computationally intensive, the IBPF algorithm enables us to directly estimate the parameters of models describing high-dimensional partially-observed nonlinear dynamic systems via likelihood maximization. The ability to directly estimate parameters of Model 3 is responsible for the large increase in model likelihoods reported in Table 1, as no algorithm with similar capabilities existed when Lee et al. (2020a) published their results. Simulations from the fitted model are displayed in Fig. 4.

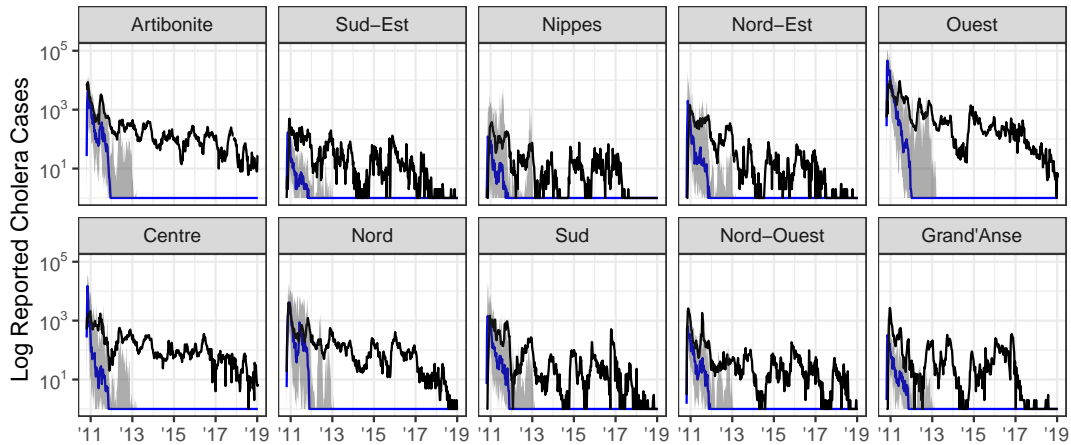


FIG 4. Simulations from initial conditions using the spatially coupled version of Model 3. The black curve represents true case count, the blue line the median of 500 simulations from the model, and the gray ribbons representing 95% confidence interval.

3.2. Model Diagnostics. Parameter calibration (whether Bayesian or frequentist) aims to find the best description of the observed data under the assumptions of the model. Obtaining the best fitting set of parameters for a given model does not, however, guarantee that the model provides an accurate representation of the system in question. Model misspecification, which may be thought of as the omission of a mechanism in the model that is an important feature of the dynamic system, is inevitable at all levels

of model complexity. To make progress, while accepting proper limitations, one must bear in mind the much-quoted observation of Box (1979) that “all models are wrong but some are useful.” Beyond being good practical advice for applied statistics, this assertion is relevant for the philosophical justification of statistical inference as severe testing (Mayo, 2018, Sec. 4.8). In this section, we provide some tools and suggestions for diagnosing mechanistic models with the goal of making the subjective assessment of model “usefulness” more objective. To do this, we will rely on the quantitative and statistical ability of the model to match the observed data, which we call the model’s *goodness-of-fit*, with the guiding principle that a model which cannot adequately describe observed data may not be reliable for useful purposes. Goodness-of-fit may provide evidence supporting the causal interpretation of one model versus another, but cannot by itself rule out the possibility of alternative explanations.

One common approach to assess a mechanistic model’s goodness-of-fit is to compare simulations from the fitted model to the observed data. Visual inspection may indicate defects in the model, or may suggest that the observed data are a plausible realization of the fitted model. While visual comparisons can be informative, they provide only a weak and informal measure of the goodness-of-fit of a model. The study by Lee et al. (2020a) provides an example of this: their models and parameter estimates resulted in simulations that visually resembled the observed data, yet resulted in model likelihoods that were—in some cases—considerably smaller than likelihoods that can be achieved via the likelihood based optimization techniques that were used (see Table 1). Alternative forms of model validation should therefore be used in conjunction with visual comparisons of simulations to observed data.

Another approach is to compare a quantitative measure of the model fit (such as MSE, predictive accuracy, or model likelihood) among all proposed models. These comparisons provide insight into how each model performs relative to the others. To calibrate relative measures of fit, it is useful to compare against a model that has well-understood statistical ability to fit data, and we call this model a *benchmark*. Standard statistical models, interpreted as associative models without requiring any mechanistic interpretation of their parameters, provide suitable benchmarks. Examples include linear regression, auto-regressive moving average time series models, or even independent and identically distributed measurements. The benchmarks enable us to evaluate the goodness of fit that can be expected of a suitable mechanistic model.

Goodness-of-fit alone does not guarantee that a model provides a correct causal interpretation of the model. Indeed, associative models are not constrained to have a causal interpretation, and typically are designed with the sole goal of providing a statistical fit to data. Consequently, we should not require a candidate mechanistic model to beat all benchmarks. However, a mechanistic model which falls far short against benchmarks is evidently failing to explain some substantial aspect of the data. A convenient measure of fit should have interpretable differences that help to operationalize the meaning of far short. Ideally, the measure should also have favorable theoretical properties. Consequently, we focus on log-likelihood as a measure of goodness of fit, and we adjust for the degrees of freedom of the models to be compared by using the Akaike information criterion (AIC) (Akaike, 1974).

In some cases, a possible benchmark model could be a generally accepted mechanistic model, but often no such model is available. Because of this, we use a log-linear Gaussian ARMA model as an associative benchmark, as recommended by He, Ionides and King (2010). The theory and practice of ARMA models is well developed, and these linear models are appropriate on a log scale due to the exponential growth and decay characteristic of biological dynamics. Likelihoods of Models 1–3 and their respective ARMA benchmark models are provided in Table 1.

It should be universal practice to present measures of goodness of fit for published models, and mechanistic models should be compared against benchmarks. In our literature review of the Haiti cholera epidemic, no benchmark models were considered in any of the 24 papers which calibrated a mechanistic model to data in order to obtain scientific conclusions. Including benchmarks would help authors and readers to detect and confront any major statistical limitations of the proposed mechanistic models. In addition, the published goodness of fit provides a concrete point of comparison for subsequent scientific investigations. When combined with online availability of data and code, objective measures of fit provide a powerful tool to accelerate scientific progress, following the paradigm of the *common task framework* (Donoho, 2017, Sec. 6).

The use of benchmarks may also be beneficial when developing models at differing spatial scales, where a direct comparison between models likelihoods is meaningless. In such a case, a benchmark model can be fit to each spatial resolution being considered, and each model compared to their respective benchmark. Large advantages (shortcomings) in model likelihood relative to the benchmark for a given spatial scale that are not present in other spatial scales may provide weak evidence for (against) the statistical fit of models across a range of spatial resolutions.

	Model 1	Model 2	Model 3
Log-likelihood	−2731.3 (−3050.1) ¹	−21957.3 (−29350.0)	−17883.2 (−42964.1) ²
Number of Fit Parameters	13 (20)	16 (26)	29 (29)
AIC	5488.7 (6140.1) ¹	43946.5 (58752.1)	35824.3 (85986.2) ²
Log-ARMA(2,1) Log-likelihood	−2802.6	−18061.9	−18061.9

TABLE 1

Log-likelihood values for each models compared to their ARMA benchmarks. Values in parenthesis are corresponding values using Lee et al. (2020a) parameter estimates. ¹*The reported likelihood is an upper bound of the likelihood of the Lee et al. (2020a) model.* ²*Lee et al. (2020a) fit Model 3 to a subset of the data (March 2014 onward, excluding data from Ouest in 2015-2016). On this subset, their model has a likelihood of −9003.0. On this same subset, our model has a likelihood of −7367.5. See Sec. S5 of the supplement material for more details on estimating the likelihood of the Lee et al. (2020a) models.*

Similar to comparing log-likelihoods across models, an additional powerful diagnosis tool is the comparison of conditional log-likelihoods. Conditional likelihoods, defined as the density $f_{Y_k|Y_1, \dots, Y_{k-1}}(Y_k = y_k^* | y_{1:k-1}^*)$, provide a basic description of how well the proposed model can describe each data point, given the previous observations. Comparing these results across models—including benchmark models—can help researchers identify potential model deficiencies, or errors in the observed data. Additional tools for assessing the goodness-of-fit of a model include plotting the effective sample size of each observation (Liu, 2001), and comparing summary statistics of the observed data to simulations from the model (Wood, 2010). These summary statistics are sometimes called diagnostic probes (King et al., 2015; King, Nguyen and Ionides, 2016).

3.3. Forecasts. Forecasts aim to provide an accurate estimate of the future state of a system based on currently available data, together with an assessment of uncertainty.

Mechanistic forecasting models, compatible with current scientific understanding, may also provide estimates of the future effects of potential interventions. Further, they may enable real-time testing of new scientific hypotheses (Lewis et al., 2022).

Recent information about a dynamic system should be more relevant for a forecast than older information. This assertion may seem self-evident, but it is not the case for deterministic models, for which the initial conditions together with the parameters are sufficient for forecasting, and so recent data do not have special importance. Epidemiological forecasts based on deterministic models are not uncommon in practice, despite their limitations (King et al., 2015). That may explain why Lee et al. (2020a) chose to obtain forecasts by simulating the calibrated models forward from initial conditions. Here, we compare with a forecast projected from the filtered distribution: to obtain a forecast using POMP model m at a collection of times $t_{N+1:N+s}$, where N is the index for the last available data, we simulate forward from a draw from $f_{\mathbf{x}_N|\mathbf{y}_{1:N}}^{(m)}(\mathbf{x}_N|\mathbf{y}_{1:N}^*; \hat{\theta})$ where $\hat{\theta}$ is a vector of calibrated parameters. The decision not to do this partially explains the unsuccessful forecasts of Lee et al. (2020a): their Table S7 shows that the subset of their simulations which were consistent with observing zero cases in 2019 predicted the elimination of cholera. In this case study, projecting the future state of the cholera epidemic starting from latent states that are draws from the filtering distribution allows the model-based forecasts to benefit from the fact that very few cholera cases were observed in 2018 and January 2019.

Uncertainty in just a single parameter can lead to drastically different forecasts (Saltelli et al., 2020). Therefore, parameter uncertainty should also be considered when obtaining model forecasts to influence policy. If a Bayesian technique is used for parameter estimation, a natural way to account for parameter uncertainty is to obtain simulations from the model where each simulation is obtained using parameters drawn from the estimated posterior distribution. For frequentist inference, an empirical Bayes approach leads to similar methodology (King et al., 2015); details are provided in the supplement. Both of these approaches share the similarity that parameters are chosen for the forecast approximately in proportion to their corresponding value of the likelihood function, $f_{\mathbf{y}_{1:N}}^{(m)}(\mathbf{y}_{1:N}^*; \theta)$.

The primary forecasting goal of Lee et al. (2020a) was to investigate the potential consequences of vaccination interventions on a system to inform policy. Outcomes of their study include estimates for the probability of cholera elimination and cumulative number of cholera infections under several possible vaccination scenarios. Mimicking their efforts, we define cholera elimination as having less than one infection of cholera over at least 52 consecutive weeks, and provide forecasts under the following vaccination scenarios:

- V0: No additional vaccines are administered.
- V1: Vaccination limited to the departments of Centre and Artibonite, deployed over a two-year period.
- V2: Vaccination limited to three departments: Artibonite, Centre, and Ouest deployed over a two-year period.
- V3: Countrywide vaccination implemented over a five-year period.
- V4: Countrywide vaccination implemented over a two-year period.

Simulations from probabilistic models (Models 1 and 3) represent possible trajectories of the dynamic system under the scientific assumptions of the models. Estimates of the probability of cholera elimination can therefore be obtained as the proportion of simulations from these models that result in cholera elimination. The results of these projections are summarized in Figs. 5–8, and suggest that cholera elimination was likely

even without increased vaccination efforts—consistent with observed reality (Trevisin et al., 2022).

Probability of elimination estimates of this form are not meaningful for deterministic models, as the trajectory of these models only represent the mean behavior of the system rather than individual potential outcomes. We therefore do not provide probability of elimination estimates under Model 2. Still, trajectories obtained by Model 2 are consistent with the simulation results of Models 1 and 3, and suggest that cholera was in the process of being eliminated from Haiti.

In addition to probability of elimination estimates, we provide estimates for the cumulative number of infections under each vaccination scenario from February 2019 – February 2024. Notably, the median number of cumulative cholera infections under the no-vaccination scenario using Models 1 and 3 were 3,486 and 1,010, respectively. While there is remaining time during this projection period in which new cholera infections can be detected, up to this point our estimates are far more consistent with the observed number of reported cholera cases than the corresponding estimates of Lee et al. (2020a), which were approximately 400,000 and 1,000,000.

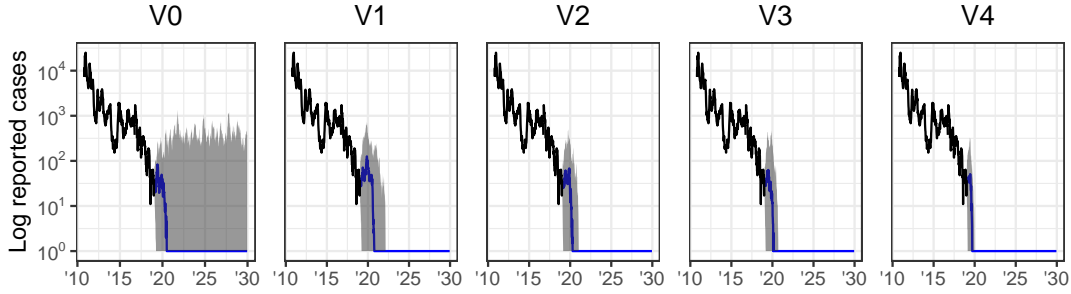


FIG 5. Simulations of Model 1 under each vaccination scenario. Blue line indicates the simulated median of reported cases, and the ribbon represents 95% of simulations.

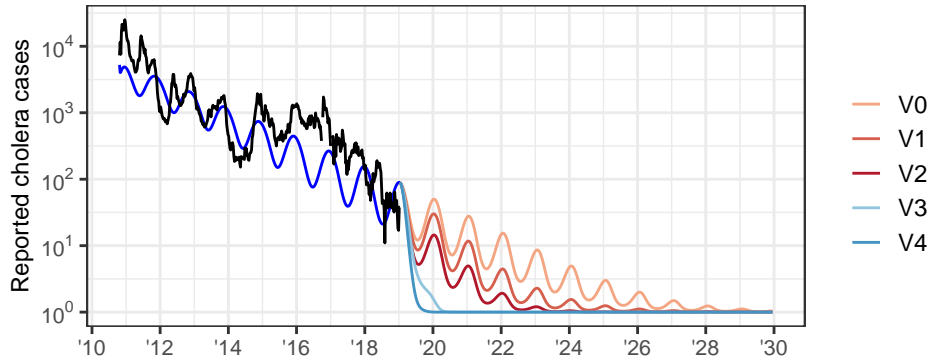


FIG 6. Simulated trajectory of Model 2 (blue curve) and projections under the various vaccination scenarios. Reported cholera incidence is shown in black.

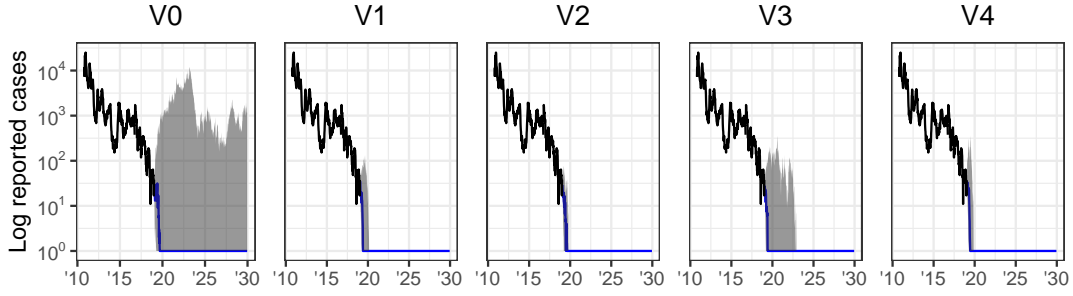


FIG 7. Simulations of Model 3 under each vaccination scenario. Blue line indicates the simulated median of reported cases, and ribbon represents 95% of simulations.

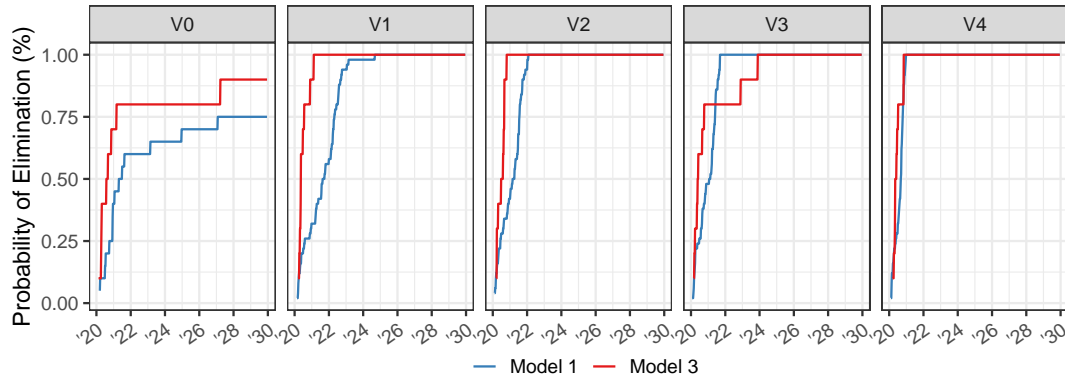


FIG 8. Probability of elimination across simulations for a 10 year period. Compare to Figure 3A of Lee et al. (2020a).

3.4. Corroborating Fitted Models with Previous Scientific Knowledge. The resulting mechanisms in a fitted model can be compared to current scientific knowledge about a system. Agreement between model based inference and our current understanding of a system may taken as a confirmation of both model based conclusions and our scientific understanding. On the other hand, comparisons may generate unexpected results that have the potential to spark new scientific knowledge (Ganusov, 2016).

In the context of our case study, we demonstrate how the fit of Model 1 corroborates other evidence concerning the role of rainfall in cholera epidemics. Specifically, we examine the results of fitting the flexible cubic spline term in Model 1 (Eq. (1)–(2)). The cubic splines permit flexible estimation of seasonality in the force of infection, $\beta(t)$. Fig. 9 shows that the estimated seasonal transmission rate β mimics the rainfall dynamics in Haiti, despite Model 1 not having access to rainfall data. This is consistent with previous studies finding that rainfall played an important role in cholera transmission in Haiti (Lemaitre et al., 2019; Eisenberg et al., 2013). The estimated seasonality also features an increased transmission rate during the Fall, which was noticed at an earlier stage of the epidemic (Rinaldo et al., 2012).

For any model-based inference, it is important to recognize and assess the modeling simplifications and assumptions that were used in order to arrive at the conclusions. In epidemiological studies, for example, quantitative understanding of individual-level processes may not perfectly match model parameters that were fit to population-level

Mechanism	Model 1	Model 2	Model 3
Infection (day)	$\mu_{IR}^{-1} = 2.0$ (6)	$\mu_{IR}^{-1} = 7.0$ (14)	$\mu_{IR}^{-1} = NA$ (26)
Latency (day)	$\mu_{EI}^{-1} = 1.4$ (5)	$\mu_{EI}^{-1} = 1.3$ (13)	—
Seasonality	$\beta_{1:6} = (1.6, 1.1, 1.3, 1.1, 1.5, 0.9)$ (2)	$a = 0.4$ (11)	$a = 4.67$ $r = 0.531$ (31)
Immunity (year)	$\mu_{RS}^{-1} = 8.0$ (7)	$\mu_{RS}^{-1} = 6.7 \times 10^{13}$ (15) $\omega_1^{-1} = 1.0$ (17) $\omega_2^{-1} = 5.0$ (17)	$\mu_{RS}^{-1} = NA$ (29)
Vaccine efficacy	—	$\theta_{1:4} = (0.80, 0.76, 0.57, 0.48)$ (12)	$\eta_{ud}(t)$
Birth/death (yr)	$\mu_S^{-1} = 44.9$ $\delta^{-1} = 134.2$ (9)	—	$\delta^{-1} = NA$ (27)
Symptomatic frac.	$f_z(t) = c\theta^*(t - \tau_d)$ (4-5)	$f = 0.2$ (13)	$f = 0.22$ (25)
Asymptomatic infectivity	$\epsilon = 0.05$ (1)	$\epsilon = 0.001$ (11) $\epsilon_W = 10^{-7}$ (19)	$\epsilon = 1$ (23) $\epsilon_W = 0.262$ (31)
Human to human	$\beta_{1:6}$ as above (1)	$\beta = 4.87 \times 10^{-17}$ (11)	$\beta_{1:10} = (4.14, 2.66, 1.70, 0.12, 0.17, 2.59, 1.07, 3.02, 1.84, 0.33) \times 10^{-6}$ (23)
Water to human	—	$W_{sat} = 10^5$ (11) $\beta_W = 1.1$	$\beta_{W1:10} = (2.07, 42.59, 2.99, 23.86, 7.94, 28.84, 5.29, 8.91, 9.40, 3.48)$ (23)
Human to water	—	$\mu_W = 9280$ (19)	$\mu_W = 5.27 \times 10^{-5}$ (31)
Water survival (wk)	—	$\delta_W^{-1} = 3$ (20)	$\delta_W^{-1} = 0.16$ (32)
Mixing exponent	$\nu = 0.97$ (1)	—	—
Process noise(wk ^{1/2})	$\sigma_{proc} = (0.31, 0.35)$ (1)	—	$\sigma_{proc} = 0.094$ (25)
Reporting rate	$\rho = 0.898$ (S16)	$\rho = 0.20$ (S17)	$\rho = 0.41$ (S18)
Observation overdispersion	$\psi = (356.16, 59.55)$ (S16)		$\psi = 63.54$ (S18)

TABLE 2

References to the relevant equation are given in parentheses. Parameters in blue were fixed based on scientific reasoning and not fitted to the data. [N] denotes parameters added during our re-analysis, not considered by Lee et al. Translations back into the notation of Lee et al. (2020a) are given in Table S1.

case counts, even when the model provides a strong statistical fit (He, Ionides and King, 2010). This makes direct interpretation of estimated parameters delicate.

Our case study provides an example of this in the parameter estimate for the duration of natural immunity due to cholera infection μ_{RS}^{-1} . Under the framework of Model 2, the best estimate for this parameter is 6.7×10^{13} , suggesting that individuals have effectively permanent immunity to cholera once infected, and thus the pool of susceptible individuals decreases over time. To interpret this result, we bear in mind that the data ranged from 2010-2019, and therefore estimates of immunity longer than 10 years—the upper end of previous estimates of natural immunity (King et al., 2008)—effectively result in the same model dynamics. The depletion of susceptible individuals may also be attributed to confounding mechanisms—such as localized vaccination programs and non-pharmaceutical interventions that reduce cholera transmission (Trevisin et al., 2022; Rebaudet et al., 2021)—that were not accounted for in the model. Perhaps the best interpretation of the estimated parameter, then, is that under model assumptions, the

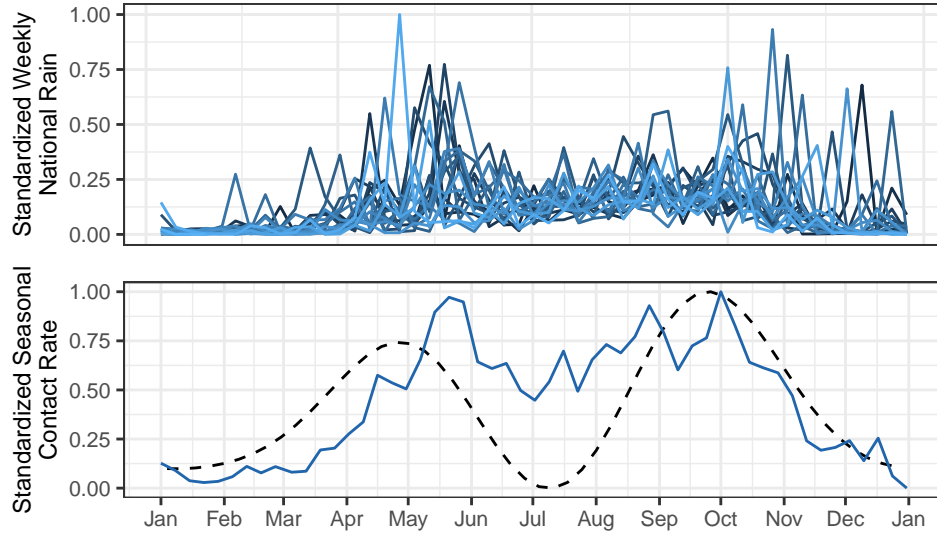


FIG 9. (Top) weekly rainfall in Haiti, lighter colors representing more recent years. (Bottom) estimated seasonality in the transmission rate (dashed line) plotted alongside mean rainfall (solid line).

model most adequately describes the observed data by having a steady decrease in the number of susceptible individuals. The weak statistical fit of Model 2 compared to a log-linear benchmark (see Table 1) cautions us against drawing quantitative conclusions from this model. A model that has a poor statistical fit may nevertheless provide a useful conceptual framework for thinking about the system under investigation. However, a claim that the model has been validated against data should be reserved for situations where the model provides a statistical fit that is competitive against alternative explanations.

4. Robust interpretation of model based conclusions. A model which aspires to provide quantitative guidance for assessing interventions should provide a quantitative statistical fit for available data. However, strong statistical fit does not guarantee a correct causal structure: it does not even necessarily require the model to assert a causal explanation. A causal interpretation is strengthened by corroborative evidence. For example, reconstructed latent variables (such as numbers of susceptible and recovered individuals) should make sense in the context of alternative measurements of these variables (Grad, Miller and Lipsitch, 2012). In addition, parameters that have been calibrated to data should make sense in the context of alternative lines of evidence about the phenomena being modeled, while making allowance for the possibility that the interpretations of parameters may vary when modeling across differing spatial scales.

If a mechanistic model including a feature (such as a representation of a mechanism, or the inclusion of a covariate) fits better than mechanistic models without that feature, and also has competitive fit compared to associative benchmarks, this may be taken as evidence supporting the scientific relevance of the feature. As for any analysis of observational data, we must be alert to the possibility of confounding. For a covariate, this shows up in a similar way to regression analysis: the covariate under investigation could be a proxy for some other unmodeled or unmeasured covariate. For a mechanism, the model feature could in principle explain the data by helping to account for some different unmodeled phenomenon. In the context of our analysis, the estimated trend

in transmission rate could be explained by any trending variable (such as hygiene improvements, or changes in population behavior), resulting in confounding from collinear covariates. Alternatively, the trend could be attributed to a decreasing reporting rate rather than decreasing transmission rate, resulting in confounded mechanisms. The robust statistical conclusion is that a model which allows for change fits better than one which does not—we argue that a decreasing transmission rate is a plausible way to explain this, but the incidence data themselves do not provide enough information to pin down the mechanism.

5. Discussion. The ongoing global COVID-19 pandemic has provided a clear example on how government policy may be affected by the conclusions of scientific models (Saltelli et al., 2020). This article demonstrates that fitting appropriate scientific models remains a challenging statistical task, and therefore great care is needed when fitting scientific models for policy recommendations. We provided a few suggestions that may aid the fitting of mechanistic models such as comparing model likelihoods to a benchmark. Improved model fits allows for meaningful statistical inference that may provide valuable insight on a dynamic system and may improve the accuracy of model forecasts. Caution is nonetheless needed when making policy based on modeling conclusions, as model misspecification may invalidate conclusions.

In this article we argue that careful attention to important statistical details could have led the models proposed by Lee et al. (2020a) to correctly predict the imminent cholera elimination. We acknowledge the benefit of hindsight: our demonstration of a statistically principled route to obtain better-fitting models with more accurate predictions does not rule out the possibility of discovering other models that fit well yet predict poorly. We used the same data and models, and even much of the same code, as Lee et al. (2020a), and yet ended up with drastically different conclusions. At a minimum, we have shown that the conclusions are sensitive to details in how the data analysis is carried out, and that attention to statistical fit (including numerical issues such as likelihood maximization) can lead to improved policy guidance.

Inference for mechanistic time series models offers opportunities for understanding and controlling complex dynamic systems. This case study has investigated issues requiring attention when applying powerful new statistical techniques that can enable statistically efficient inference for a general class of partially observed Markov process models. Researchers should check that the computationally intensive numerical calculations are carried out adequately. Comparison against benchmarks and alternative model specifications should be considered to evaluate the statistical goodness-of-fit. Once that is accomplished, care is required to assess what causal conclusions can properly be inferred given the possibility of alternative explanations consistent with the data. Studies that combine model development with thoughtful data analysis, supported by a high standard of reproducibility, build knowledge about the system under investigation. Cautionary warnings about the difficulties inherent in understanding complex systems (Saltelli et al., 2020; Ioannidis, Cripps and Tanner, 2020; Ganusov, 2016) should motivate us to follow best practices in data analysis, rather than avoiding the challenge.

5.1. Reproducibility and Extendability. Lee et al. (2020a) published their code and data online, and this reproducibility facilitated our work. Robust data analysis requires not only reproducibility but also extendability: if one wishes to try new model variations, or new approaches to fitting the existing models, or plotting the results in a different way, this should be not excessively burdensome. Scientific results are only trustworthy so far as they can be critically questioned, and an extendable analysis should facilitate such examination (Gentleman and Temple Lang, 2007).

We provide a strong form of reproducibility, as well as extendability, by developing our analysis in the context of a software package, `haitipkg`, written in the R language (R Core Team, 2022). Using a software package mechanism supports documentation, standardization and portability that promote extendability. In the terminology of Gentleman and Temple Lang (2007), the source code for this article is a *dynamic document* combining code chunks with text. In addition to reproducing the article, the code can be extended to examine alternative analysis to that presented. The dynamic document, together with the R packages, form a *compendium*, defined by Gentleman and Temple Lang (2007) as a distributable and executable unit which combines data, text and auxiliary software (the latter meaning code written to run in a general-purpose, portable programming environment, which in this case is R).

Acknowledgments. The authors would like to thank Mercedes Pascual and Betz Halloran for helpful discussions. Laura Matrajt provided additional data for the Model 2 analysis.

Funding. This work was supported by National Science Foundation grants DMS-1761603 and DMS-1646108.

SUPPLEMENTARY MATERIAL

Eliminating cholera in Haiti: Supplement

This document contains additional details for Models 1–3, as well as a translation table that facilitates comparisons between these models and those described by Lee et al. (2020a). The supplement also demonstrates our capability to faithfully replicate the results of Lee et al. (2020a).

Auxiliary software: `haitipkg`

The `haitipkg` R package is maintained in a GitHub repository, `jeswheel/haitipkg`. The submitted version of this package will be archived on Zenodo. The package contains all of the data and code used to create and fit the models, as well as other useful functions that were used in this article.

Dynamic document: `haiti_article`

The Rnw (R noweb, Ramsey, 1994) files generating this article and its supplement are maintained in a GitHub repository, `jeswheel/haiti_article`. The submitted version will be archived on Zenodo.

REFERENCES

- AKAIKE, H. (1974). A New Look at the Statistical Model Identification. *IEEE Transactions on Automatic Control* **19** 716–723.
- ANDRIEU, C., DOUCET, A. and HOLENSTEIN, R. (2010). Particle Markov chain Monte Carlo methods. *Journal of the Royal Statistical Society, Series B (Statistical Methodology)* **72** 269–342.
- ARULAMPALAM, M. S., MASKELL, S., GORDON, N. and CLAPP, T. (2002). A Tutorial on Particle Filters for Online Nonlinear, Non-Gaussian Bayesian Tracking. *IEEE Transactions on Signal Processing* **50** 174 – 188.
- ASFAW, K., IONIDES, E. L. and KING, A. A. (2021). `spatPomp`: R package for Statistical Inference for Spatiotemporal Partially Observed Markov Processes. <https://github.com/kidusasfaw/spatPomp>.
- BEHREND, M. R., BASÁÑEZ, M.-G., HAMLEY, J. I. D., PORCO, T. C., STOLK, W. A., WALKER, M., DE VLAS, S. J. and FOR THE NTD MODELLING CONSORTIUM (2020). Modelling for Policy: The Five Principles of the Neglected Tropical Diseases Modelling Consortium. *PLOS Neglected Tropical Diseases* **14** 1–17.

- BOX, G. E. (1979). Robustness in the Strategy of Scientific Model Building. In *Robustness in statistics* 201–236. Elsevier.
- BRAUER, F. (2017). Mathematical Epidemiology: Past, Present, and Future. *Infectious Disease Modelling* **2** 113–127.
- BRETÓ, C. and IONIDES, E. L. (2011). Compound Markov Counting Processes and their Applications to Modeling Infinitesimally Over-Dispersed Systems. *Stochastic Processes and their Applications* **121** 2571–2591.
- BRETÓ, C., HE, D., IONIDES, E. L. and KING, A. A. (2009). Time Series Analysis via Mechanistic Models. *Annals of Applied Statistics* **3** 319–348.
- COLLINS, O. C. and GOVINDER, K. S. (2014). Incorporating Heterogeneity into the Transmission Dynamics of a Waterborne Disease Model. *Journal of Theoretical Biology* **356** 133–143.
- DADLANI, A., AFOLABI, R. O., JUNG, H., SOHRABY, K. and KIM, K. (2020). Deterministic Models in Epidemiology: From Modeling to Implementation.
- DONOHO, D. (2017). 50 Years of Data Science. *Journal of Computational and Graphical Statistics* **26** 745–766.
- EISENBERG, M. C., KUJBIDA, G., TUITE, A. R., FISMAN, D. N. and TIEN, J. H. (2013). Examining Rainfall and Cholera Dynamics in Haiti using Statistical and Dynamic Modeling Approaches. *Epidemics* **5** 197–207.
- FRANCOIS, J. (2020). Cholera Remains a Public Health Threat in Haiti. *The Lancet Global Health* **8** e984.
- GANUSOV, V. V. (2016). Strong Inference in Mathematical Modeling: a Method for Robust Science in the Twenty-First Century. *Frontiers in Microbiology* **7** 1131.
- GENTLEMAN, R. and TEMPLE LANG, D. (2007). Statistical Analyses and Reproducible Research. *Journal of Computational and Graphical Statistics* **16** 1–23.
- GIORDANO, G., BLANCHINI, F., BRUNO, R., COLANERI, P., DI FILIPPO, A., DI MATTEO, A. and COLANERI, M. (2020). Modelling the COVID-19 Epidemic and Implementation of Population-Wide Interventions in Italy. *Nature Medicine* **26** 855–860.
- GRAD, Y. H., MILLER, J. C. and LIPSITCH, M. (2012). Cholera Modeling: Challenges to Quantitative Analysis and Predicting the Impact of Interventions. *Epidemiology (Cambridge, Mass.)* **23** 523.
- GREEN, K. C. and ARMSTRONG, J. S. (2015). Simple Versus Complex Forecasting: The Evidence. *Journal of Business Research* **68** 1678–1685. Special Issue on Simple Versus Complex Forecasting.
- HE, D., IONIDES, E. L. and KING, A. A. (2010). Plug-and-Play Inference for Disease Dynamics: Measles in Large and Small Towns as a Case Study. *Journal of the Royal Society Interface* **7** 271–283.
- HENRYS, J. H., LEREBOURS, G., ACHILLE, M. A., MOISE, K. and RACCURT, C. (2020). Cholera in Haiti. *The Lancet Global Health* **8** e1469.
- IOANNIDIS, J. P., CRIPPS, S. and TANNER, M. A. (2020). Forecasting for COVID-19 has Failed. *International Journal of Forecasting*.
- IONIDES, E. L., NING, N. and WHEELER, J. (2022). An Iterated Block Particle Filter for Inference on Coupled Dynamic Systems with Shared and Unit-Specific Parameters.
- IONIDES, E. L., BRETO, C., PARK, J., SMITH, R. A. and KING, A. A. (2017). Monte Carlo Profile Confidence Intervals for Dynamic Systems. *Journal of the Royal Society Interface* **14** 1–10.
- KERMACK, W. O. and MCKENDRICK, A. G. (1927). A Contribution to the Mathematical Theory of Epidemics. *Proceedings of the Royal Society of London, Series A* **115** 700–721.
- KING, A. A., NGUYEN, D. and IONIDES, E. L. (2016). Statistical Inference for Partially Observed Markov Processes via the R Package pomp. *Journal of Statistical Software* **69** 1–43.
- KING, A. A. and ROWAN, T. (2020). subplex: Unconstrained Optimization using the Subplex Algorithm R package version 1.6.
- KING, A. A., IONIDES, E. L., PASCUAL, M. and BOUMA, M. J. (2008). Inapparent Infections and Cholera Dynamics. *Nature* **454** 877–880.
- KING, A. A., IONIDES, E. L., BRETÓ, C. M., ELLNER, S. and KENDALL, B. (2009). pomp: Statistical Inference for Partially Observed Markov Processes. R package, available at <http://cran.r-project.org/web/packages/pomp>.
- KING, A. A., DOMENECH DE CELLÈS, M., MAGPANTAY, F. M. and ROHANI, P. (2015). Avoidable Errors in the Modelling of Outbreaks of Emerging Pathogens, with Special Reference to Ebola. *Proceedings of the Royal Society B: Biological Sciences* **282** 20150347.
- LEE, E. C., CHAO, D. L., LEMAITRE, J. C., MATRAJT, L., PASETTO, D., PEREZ-SAEZ, J., FINGER, F., RINALDO, A., SUGIMOTO, J. D., HALLORAN, M. E., LONGINI, I. M., TERNIER, R., VISSIERES, K., AZMAN, A. S., LESSLER, J. and IVERS, L. C. (2020a). Achieving Coordinated

- National Immunity and Cholera Elimination in Haiti Through Vaccination: A Modelling Study. *The Lancet Global Health* **8** e1081–e1089.
- LEE, E. C., TERNIER, R., LESSLER, J., AZMAN, A. S. and IVERS, L. C. (2020b). Cholera in Haiti—Authors’ Reply. *The Lancet Global Health* **8** e1470–e1471.
- LEMAITRE, J., PASETTO, D., PEREZ-SAEZ, J., SCIARRA, C., WAMALA, J. F. and RINALDO, A. (2019). Rainfall as a Driver of Epidemic Cholera: Comparative Model Assessments of the Effect of Intra-Seasonal Precipitation Events. *Acta Tropica* **190** 235–243.
- LEWIS, A. S. L., ROLLINSON, C. R., ALLYN, A. J., ASHANDER, J., BRODIE, S., BROOKSON, C. B., COLLINS, E., DIETZE, M. C., GALLINAT, A. S., JUVIGNY-KHENAFOU, N., KOREN, G., MCGLINN, D. J., MOUSTAHFID, H., PETERS, J. A., RECORD, N. R., ROBBINS, C. J., TONKIN, J. and WARDLE, G. M. (2022). The Power of Forecasts to Advance Ecological Theory. *Methods in Ecology and Evolution*.
- LIU, J. S. (2001). *Monte Carlo Strategies in Scientific Computing*. Springer, New York.
- LUCAS, R. E. et al. (1976). Econometric Policy Evaluation: A Critique. In *Carnegie-Rochester conference series on public policy* **1** 19–46.
- MAYO, D. G. (2018). *Statistical Inference as Severe Testing*. Cambridge: Cambridge University Press.
- NDII, M. Z. and SUPRIATNA, A. K. (2017). Stochastic Mathematical Models in Epidemiology. *Information* **20** 6185–6196.
- NING, N. and IONIDES, E. L. (2021). Iterated Block Particle Filter for High-dimensional Parameter Learning: Beating the Curse of Dimensionality.
- PARK, J. and IONIDES, E. L. (2020). Inference on High-Dimensional Implicit Dynamic Models using a Guided Intermediate Resampling Filter. *Statistics & Computing* **30** 1497–1522.
- PEZZOLI, L. (2020). Global Oral Cholera Vaccine Use, 2013–2018. *Vaccine* **38** A132–A140. Cholera Control in Three Continents: Vaccines, Antibiotics and WASH.
- RAMSEY, N. (1994). Literate Programming Simplified. *IEEE software* **11** 97–105.
- REBAUDET, S., GAUDART, J. and PIARROUX, R. (2020). Cholera in Haiti. *The Lancet Global Health* **8** e1468.
- REBAUDET, S., BULIT, G., GAUDART, J., MICHEL, E., GAZIN, P., EVERS, C., BEAULIEU, S., ABEDI, A. A., OSEI, L., BARRAIS, R. et al. (2019). The Case-Area Targeted Rapid Response Strategy to Control Cholera in Haiti: A Four-Year Implementation Study. *PLoS Neglected Tropical Diseases* **13** e0007263.
- REBAUDET, S., DÉLY, P., BONCY, J., HENRYS, J. H. and PIARROUX, R. (2021). Toward Cholera Elimination, Haiti. *Emerging Infectious Diseases* **27** 2932.
- REBESCHINI, P. and VAN HANDEL, R. (2015). Can Local Particle Filters Beat the Curse of Dimensionality? *The Annals of Applied Probability* **25** 2809–2866.
- RINALDO, A., BERTUZZO, E., MARI, L., RIGHETTO, L., BLOKESCH, M., GATTO, M., CASAGRANDE, R., MURRAY, M., VESENBECKH, S. M. and RODRIGUEZ-ITURBE, I. (2012). Re-assessment of the 2010–2011 Haiti Cholera Outbreak and Rainfall-Driven Multiseason Projections. *Proceedings of the National Academy of Sciences* **109** 6602–6607.
- SALTELLI, A., BAMMER, G., BRUNO, I., CHARTERS, E., DI FIORE, M., DIDIER, E., NELSON ESPELAND, W., KAY, J., LO PIANO, S., MAYO, D. et al. (2020). Five Ways to Ensure that Models Serve Society: a Manifesto.
- STOCKS, T., BRITTON, T. and HÖHLE, M. (2020). Model Selection and Parameter Estimation for Dynamic Epidemic Models via Iterated Filtering: Application to Rotavirus in Germany. *Biostatistics* **21** 400–416.
- R CORE TEAM (2022). R: A Language and Environment for Statistical Computing R Foundation for Statistical Computing, Vienna, Austria.
- TONI, T., WELCH, D., STRELKOWA, N., IPSEN, A. and STUMPF, M. P. H. (2009). Approximate Bayesian Computation Scheme for Parameter Inference and Model Selection in Dynamical Systems. *Journal of the Royal Society Interface* **6** 187–202.
- TRACY, M., CERDÁ, M. and KEYES, K. M. (2018). Agent-Based Modeling in Public Health: Current Applications and Future Directions. *Annual Review of Public Health* **39** 77–94. PMID: 29328870.
- TREVISIN, C., LEMAITRE, J. C., MARI, L., PASETTO, D., GATTO, M. and RINALDO, A. (2022). Epidemicity of Cholera Spread and the Fate of Infection Control Measures. *Journal of the Royal Society Interface* **19** 20210844.
- WOOD, S. N. (2010). Statistical Inference for Noisy Nonlinear Ecological Dynamic Systems. *Nature* **466** 1102–1104.