

TUTORIAL IN BIOSTATISTICS OPEN ACCESS

The Mathematics of Serocatalytic Models With Applications to Public Health Data

Everlyn Kamau¹  | Junjie Chen² | Sumali Bajaj³ | Nicolás Torres⁴ | Richard Creswell⁵ | Jaime A. Pavlich-Mariscal⁶  | Christl Donnelly² | Zulma Cucunubá⁴ | Ben Lambert²

¹Francis I. Proctor Foundation, University of California San Francisco, San Francisco, California, USA | ²Department of Statistics & Pandemic Sciences Institute, University of Oxford, Oxford, UK | ³Department of Biology & Department of Statistics & Merton College & Pandemic Sciences Institute, University of Oxford, Oxford, UK | ⁴Instituto de Salud Pública, Pontificia Universidad Javeriana, Bogota, Colombia | ⁵Melbourne School of Population and Global Health, University of Melbourne, Melbourne, Australia | ⁶Facultad de Ingeniería, Pontificia Universidad Javeriana, Bogota, Colombia

Correspondence: Everlyn Kamau (everlyn.kamau@ucsf.edu)

Received: 24 January 2025 | **Revised:** 6 June 2025 | **Accepted:** 24 June 2025

Funding: JC is supported by the Moh Family Foundation on an Oxford-Moh Family Foundation Global Health Scholarship. SB was supported by the Clarendon Scholarship, St. Edmund Hall College, and NERC DTP [grant number NE/S007474/1], University of Oxford. This work was also supported by the UK NIHR Health Protection Research Unit (HPRU) in Emerging and Zoonotic Infections, a partnership between UKHSA, University of Oxford, University of Liverpool, and Liverpool School of Tropical Medicine (grant number NIHR200907 supporting C.A.D.). NT and ZC were supported by the TRACE-LAC project (Enhancing Tools for Response, Analytics and Control of Epidemics in Latin America and the Caribbean) [Grant No. 109848-002], funded by the International Development Research Center (IDRC).

ABSTRACT

Serocatalytic models are powerful tools which can be used to infer historical infection patterns from age-structured serological surveys. These surveys are especially useful when disease surveillance is limited and have an important role to play in providing a ground truth gauge of infection burden. In this tutorial, we consider a wide range of serocatalytic models to generate epidemiological insights. With mathematical analysis, we explore the properties and intuition behind these models and include applications to real data for a range of pathogens and epidemiological scenarios. We also include practical steps and code in R and Stan for interested learners to build experience with this modeling framework. Our work highlights the usefulness of serocatalytic models and shows that accounting for the epidemiological context is crucial when using these models to understand infectious disease epidemiology.

1 | Information Contained in Age-Structured Serological Data

In a serological survey, specimens from a sample population are tested for antibodies to a pathogen, yielding quantitative measurements that are typically categorized into binary outcomes. An individual is then *seropositive* if their antibody quantity is above a threshold and *seronegative* otherwise [1]. Serological surveys (or *serosurveys*) are often cross-sectional, performed

at a single time point. From these surveys, we calculate the proportion of seropositive individuals, that is, *seroprevalence* or alternatively *seropositivity*. We use these terms interchangeably throughout this paper. In Figure 1, we show data from three yellow fever serosurveys in the Americas as presented in Hugo Muench's classic 1934 study [2], which popularized serocatalytic modeling. Whereas the two Brazilian datasets generally display an upward trend in seroprevalence with age, the Colombian seroprevalence is generally flat across ages.

Everlyn Kamau, Junjie Chen, Sumali Bajaj and Nicolás Torres contributed equally to this work.

This is an open access article under the terms of the [Creative Commons Attribution](https://creativecommons.org/licenses/by/4.0/) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

© 2025 The Author(s). *Statistics in Medicine* published by John Wiley & Sons Ltd.

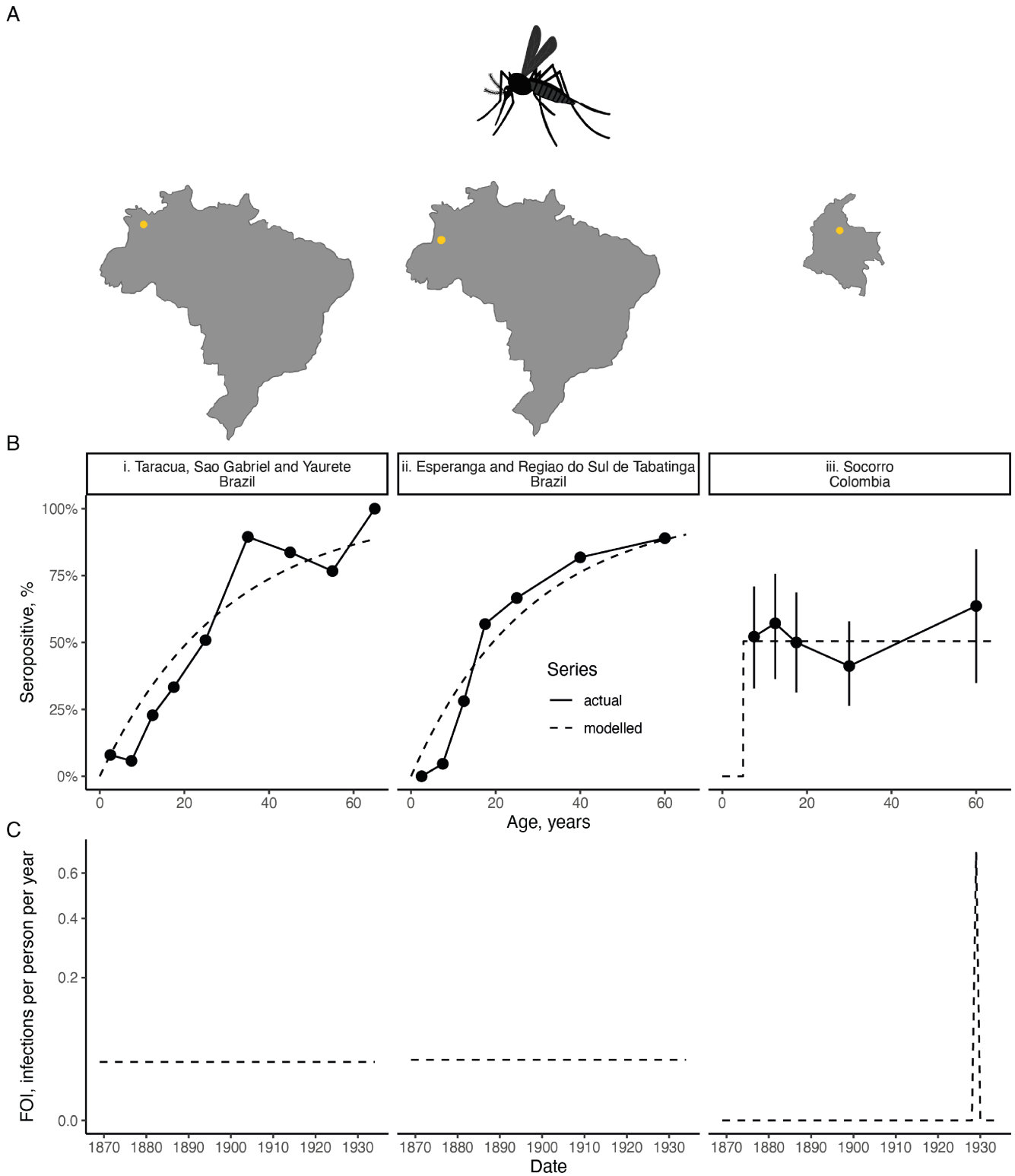


FIGURE 1 | Muench's serosurvey data for yellow fever from three locations in the Americas. Row A shows the locations of the serosurveys. Each column in Rows B and C represents a particular location in the Americas. Row B shows the raw data (black points connected by solid lines) and modeled proportions seropositive (dashed lines); in Column iii, we show the 2.5th and 97.5th percentiles of a posterior distribution assuming a uniform prior over the proportion seropositive and a binomial likelihood for each age group separately; in Columns i & ii, we do not show uncertainty intervals since we did not have access to the sample sizes used in the serosurveys. Row C shows the inferred historical FOIs. Details of the data and analysis are provided in Section A.

TABLE 1 | Summary of the serocatalytic models introduced in this study.

Model/Features	Constant FOI?	Time-dependent FOI?	Age-dependent FOI?	Sero-reversion?	Parameters	Section in Paper
Constant FOI	x			x	$\bar{\lambda}, \mu$	4
Time-dependent FOI		x		x	$\bar{\lambda}_t, \mu$	5
Age-dependent FOI			x	x	$\bar{\lambda}_a, \mu$	6
Elevated death rate due to infection	x		x		$\bar{\lambda}, \bar{\lambda}_a, \epsilon, \rho$	7
Time- and age-dependent FOI		x	x	x	$\bar{\lambda}_t, \bar{\lambda}_a, \mu$	8
Maternal antibodies	x		x		$\bar{\lambda}, \bar{\lambda}_a, \gamma$	9

Note: The models explored here are listed in the first column, and the corresponding features available with each model derivation are marked with “x” in the respective cells. The images shown in the three middle columns show the recurring pictorial representations that we use to represent these classes of models. Also shown are notations for parameters that can be inferred in each model. The section describing each model is listed in the last column. FOI: Force of infection.

With these serosurveys, Muench asks the question: what transmission patterns generated these data? He argues that the two Brazilian datasets (Figure 1B (i, ii)) were generated from a transmission history distinct from that of Colombia. Running contrary to popular opinion at the time, his claim was that yellow fever was effectively *endemic* in these locations in Brazil. In Sorocco, Colombia, however, it was known that the town had, in living memory, been free from yellow fever until an isolated outbreak occurred in 1929; that is, it was known that transmission of yellow fever was effectively *epidemic* in this region.

Transmission rate is quantified through a metric known as the *force of infection*, or FOI, the rate at which individuals become infected per unit time. FOI has units of the average number of infections a person receives per unit time (usually per year)—note, this metric can exceed 1 if it is expected that an individual would be infected more than once in a given time interval. Endemic transmission then translates into assuming that the FOI varies little over time. This assumption results in a characteristic increase in seroprevalence as individuals age, and in Figure 1B (i, ii), we show the estimated seroprevalence resulting from assuming an FOI that does not vary over time. The estimated equivalents (dashed lines) generally provide a reasonable approximation to the underlying trends in the raw data (solid points and lines).

For an isolated epidemic, assuming everyone sampled was exposed to the same degree results in a constant proportion seropositive for those old enough to have experienced and survived the epidemic. We show the estimated seropositive proportion under this assumption in Figure 1B (iii). In Figure 1C, we then show the estimated FOIs. Under these assumptions, the historical FOIs in the two Brazilian locations were virtually identical, and the inferred FOI for the 1929 Colombian yellow fever outbreak is far in excess of these.

This example illustrates the power of serological data: by making assumptions about the nature of transmission dynamics, we can reveal unobserved historical transmission patterns. These patterns can be used to project future dynamics and inform vaccination strategies [3, 4]. Serological data are critical for estimating disease burden and herd immunity thresholds, and are mostly useful for subclinical diseases that are under-recognized or under-notified [5]. Seroprevalence surveys also contribute to better analysis and interpretation of clinical surveillance data [6].

What serological data tell us about transmission patterns hinges on assumptions about how, if at all, transmission varies over time. It also depends on assumptions about exposure patterns across different demographics, and whether individuals lose antibody detectability over time—that is, *seroreversion*. Such assumptions are embodied by mathematical models or *serocatalytic* models, and this article provides an introductory guide to these models.

2 | Reproducing Our Methods and Results

This paper forms an accessible guide for serocatalytic modeling and promotes learning through step-by-step mathematical derivations, simulation of models, and fitting these to data to infer transmission dynamics. Table 1 summarizes the models covered and the structure of this article. A summary of the mathematical symbols used throughout is provided in Glossary S1. A key notational rule we use throughout is to use lowercase symbols to denote continuous quantities (e.g., b , the specific calendar date/time at which birth occurs) and uppercase symbols to denote integer values (e.g., B , the year in which individuals were born).

We demonstrate the behavior of each model by simulating from it using R [7]. The simulated dynamics are shown in figures throughout the paper, and we provide a website which includes R code and further visualizations (<https://arianajunjie.github.io/seropackage/>). We also illustrate the application of serocatalytic models by fitting actual serological data of various pathogens, with code available from https://github.com/ekamau/serocatalytic_models. For the model fitting, the parameters of each model were estimated through a Bayesian framework. More detailed information on the inference is provided in Section A. For our model fitting, we used the *targets* [8] R package to create reproducible data analysis pipelines and *renv* [9] to allow our computational environment to be reproducible by others.

3 | Transmission Dynamics Models and Their Relationship With Serocatalytic Models

Transmission dynamics models attempt to mechanistically model the population dynamics of pathogens through host populations. Many such models are built upon the classic Susceptible-Infected-Recovered (SIR) model, with the following form:

$$\begin{aligned}\frac{dS(t)}{dt} &= -\beta S(t)I(t), \\ \frac{dI(t)}{dt} &= \beta S(t)I(t) - \gamma I(t), \\ \frac{dR(t)}{dt} &= \gamma I(t)\end{aligned}\tag{1}$$

where $S(0) = S_0$, $I(0) = I_0$, $R(0) = R_0$ are the initial conditions. We assume $S(t) + I(t) + R(t) = 1$, for $t \geq 0$, meaning we model the dynamics of the proportion of individuals who are susceptible, ($S(t)$), infected ($I(t)$) and recovered (meaning they cannot be reinfected; $R(t)$).

Models such as Equation (1) attempt to mechanistically describe how an epidemic system unfolds with time. But this comes at a cost—to predict future system behaviors, it is necessary to have reasonable knowledge about the system. To solve Equation (1), we need to know the transmission rate, $\beta > 0$, the rate of recovery from infection, $\gamma > 0$, and the initial conditions. For novel pathogens or during an outbreak, these parameter values may be unknown.

Further, SIR models like Equation (1) are highly idealized, and this model system does not normally allow β or γ to vary over

time; for example, in response to the imposition of public health interventions or due to the introduction of novel pathogen variants. Equation (1) is also inherently a short-term model system since it fails to account for births or natural deaths, which collectively change the susceptible population. The model above also does not include seroreversion or account for loss of antibody detection, where individuals once recovered become susceptible or *seronegative*. For a review of the various ways in which simple SIR-type models may be extended to answer a variety of epidemiological questions, see [10].

Transmission dynamics models have been extraordinarily successful, but their correct use requires incorporating substantial epidemiological knowledge. Serocatalytic models are limited. Whereas some classes of transmission dynamics models can be used to predict future incidence, serocatalytic models cannot be used directly to predict how an epidemic might unfold into the future—they inherently look into the past.

Whereas transmission dynamics models track the status of a population through time, including the proportions infected (Figure 2A), serocatalytic models instead attempt to explain and quantify the rates at which individuals have historically become infected throughout their lifetimes (Figure 2C). Whereas

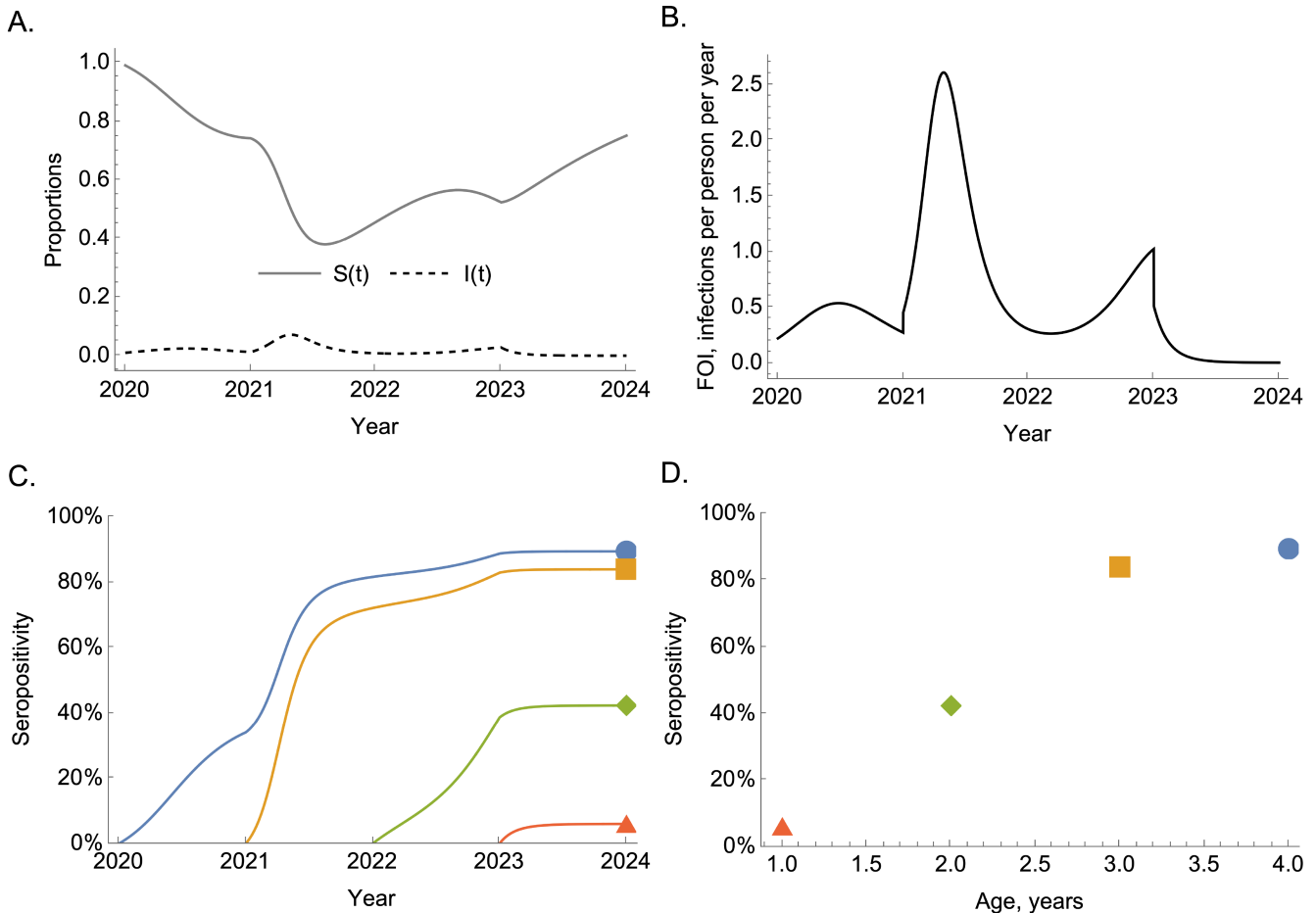


FIGURE 2 | The relationship between transmission dynamics models and serocatalytic models. Panel A shows the proportions susceptible and infected resulting from simulating a transmission dynamics model that includes waning immunity. In this model, the transmission rate, β , varies over time (see Section A for a complete description). Panel B shows the force of infection: $\beta(t)I(t)$. Panel C shows the seropositivity trajectories of four birth cohorts. Panel D shows the serological age profile, in years, in the population in 2024.

transmission dynamics models provide a mechanistic basis for *how* infected individuals contribute to ongoing transmission, serocatalytic models do not directly model the drivers of ongoing transmission but instead are useful to infer *when* infections occurred.

We can convert the SIR example of Equation (1) into a serocatalytic model by effectively ignoring the relatively short-term dynamics of the infected population and considering infection dynamics within just two groups or compartments: (i) susceptible (or seronegative), $S(t)$, and (ii) *seropositive*, $X(t)$ —those previously infected, which includes both currently infected and recovered. To remove dependence on the numbers of infected individuals and to generalize the system to allow for time-varying transmission rates, we introduce $\lambda(t) = \beta(t)I(t)$, where $\lambda(t)$ is the FOI (Figure 2B). Equation (1) then simplifies to:

$$\begin{aligned}\frac{dS(t)}{dt} &= -\lambda(t)S(t), \\ \frac{dX(t)}{dt} &= \lambda(t)S(t)\end{aligned}\quad (2)$$

However, Equation (2) requires knowing an initial condition: the proportion susceptible at some time in the past, represented by $0 \leq S(0) \leq 1$, which may not be known if no serosurveys were previously conducted. For this, serocatalytic models both simplify and complicate the problem simultaneously—they abandon the idea of modeling the entire population and instead consider birth cohorts, individuals born in the same year. Whilst this means that we get one set of equations for each birth cohort, we gain an initial condition by assuming that each birth cohort is seronegative at birth (see¹). This means that, unlike transmission dynamics models, serocatalytic models are *always* age-structured. We assume that each birth cohort has the same birth time, $b \leq t$.

For each cohort, then, we have equations of the form:

$$\begin{aligned}\frac{dS^b(t)}{dt} &= -\lambda(t)S^b(t), \\ \frac{dX^b(t)}{dt} &= \lambda(t)S^b(t)\end{aligned}\quad (3)$$

where $0 \leq S^b(t) \leq 1$ is the proportion of individuals born at time b who are seronegative at time t and have initial conditions: $S^b(b) = 1$, $X^b(b) = 0$ meaning that $S^b(t) + X^b(t) = 1$ for all $t \geq b$.

Using $S^b(t) = 1 - X^b(t)$, we can rewrite the second equation for those seropositive in Equation (3) as:

$$\frac{dX^b(t)}{dt} = \lambda(t)(1 - X^b(t))\quad (4)$$

By solving Equation (4) for each birth cohort, we can then determine their serological trajectories through time (Figure 2C); in Section 4 we explain how this equation can be solved. This allows us to determine a cross-sectional profile of seropositivity stratified by age at a particular point in time (Figure 2D); note that the age profile of seropositivity may vary for serosurveys conducted at different points in time.

3.1 | Population-Level Versus Individual-Level Quantities

FOI represents the average number of infections a person receives per unit time, and the typical time unit is years; it is a rate, meaning it must be zero or above and has no upper bound. An alternative way to represent this quantity is to convert it into the probability that an individual becomes infected in a particular year. To obtain this quantity, we consider the risk that an individual becomes infected within a year. If the FOI within a year is constant at $\bar{\lambda}$, the probability an individual remains uninfected (U) throughout that year can be calculated by solving:

$$\frac{dU(t)}{dt} = -\bar{\lambda}U(t)\quad (5)$$

with $U(0) = 1$ meaning the individual starts in an uninfected state, there is a negative sign in Equation (5) because individuals flow from being uninfected to being infected over time. To solve this for the probability they are uninfected at the end of year, we use the method of separation of variables to yield: $U(1) = \exp(-\bar{\lambda})$ meaning the probability they are infected is $1 - \exp(-\bar{\lambda})$.

More generally, when using serocatalytic models, there is a distinction between population quantities and individual quantities. The seroprevalence is the proportion of individuals who are seropositive in the population; the analogous quantity for a randomly chosen individual is the probability that they are seropositive, and these two quantities will be numerically equal and are given by $X^b(t)$. For example, if the seroprevalence is 50%, the probability that a randomly chosen individual is seropositive is 0.5.

4 | Constant Force of Infection

We first consider an idealized situation where the FOI is invariant over time, representing *endemic* transmission thought to have resulted in the transmission patterns for yellow fever in the Brazilian Amazon, see Section 1 [2]. This type of FOI could not be generated by a transmission dynamics model like that described by Equation (1) because it permits no non-zero disease equilibrium, but a simple modification of this equation to include births of susceptible individuals and/or waning immunity would allow this [10].

Assuming the FOI is fixed at $\lambda(t) = \bar{\lambda}$, Equation (4) becomes:

$$\frac{dX^b(t)}{dt} = \bar{\lambda}(1 - X^b(t))\quad (6)$$

We can then separate the variables:

$$\frac{dX^b(t)}{dt} + \bar{\lambda}X^b(t) = \bar{\lambda}\quad (7)$$

and employ the “integrating factor” approach to solve differential equations by multiplying both sides of Equation (7) by $\exp(\bar{\lambda}t)$ and rewriting the left-hand side as a derivative:

$$\frac{d}{dt} \left(X^b(t) \exp(\bar{\lambda}t) \right) = \bar{\lambda} \exp(\bar{\lambda}t)\quad (8)$$

We then integrate both sides from b , when individuals were born, to $t > b$. To do so, we introduce a dummy integration variable t' :

$$\int_b^t \frac{d}{dt'} \left(X^b(t') \exp(\bar{\lambda}t') \right) dt' = \int_b^t \bar{\lambda} \exp(\bar{\lambda}t') dt' \quad (9)$$

which leads to the following:

$$X^b(t) \exp(\bar{\lambda}t) - X^b(b) \exp(\bar{\lambda}b) = \exp(\bar{\lambda}t) - \exp(\bar{\lambda}b) \quad (10)$$

Using the initial conditions, we can rearrange Equation (10) to obtain the proportion seropositive at time t :

$$X^b(t) = 1 - \exp(-\bar{\lambda}(t - b)) \quad (11)$$

We can use Equation (11) for a range of birth cohorts to determine their seropositivity trajectories through time (to produce seropositivity profiles like those shown in Figure 2).

4.1 | Seropositivity by Age

Another way to write Equation (11) is to introduce a birth-cohort-specific age, $a^b := t - b$, for any $t \geq b$, and b is constant:

$$X^b(a^b) = 1 - \exp(-a^b \bar{\lambda}) \quad (12)$$

which shows that, as individuals age, their seropositivity increases due to their cumulative risk of infection.

We can also consider seropositivity across a population at a snapshot at time t . For this, we introduce a variable, $a_t = t - b$, to denote the age of individuals at time t ; note that, unlike for our birth-cohort-specific age, b varies since individuals in the population may be of different ages. We can then use Equation (12) to determine how seropositivity varies across the population:

$$X_t(a_t) = 1 - \exp(-a_t \bar{\lambda}) \quad (13)$$

While Equations (12) and (13) appear similar, they represent quite different quantities: The former gives the seropositivity trajectory for a particular birth-cohort born at time b , for example, like the trajectories shown in Figure 2C; the latter gives the average seropositivity for an individual of a specific age at time t (i.e., across all birth cohorts born prior to that year)—see Figure 2D for an example of such a serological cross-section.

5 | Time-Varying Force of Infection

A more general situation is when the FOI has varied historically. It can be assumed that the FOI is a continuous function of time (e.g., through the use of splines), but, for many practical purposes, ages of individuals in serological surveys are recorded in calendar years. So, often a sensible simplifying assumption is that the FOI is piecewise-constant, typically with pieces of one year in length. In what follows, we assume this unless explicitly stated otherwise. Under the piecewise-constant assumption, the FOI experienced by a cohort born at the start of year B until year $T > B$ is given by a series of yearly FOIs: $\{\bar{\lambda}_B, \bar{\lambda}_{B+1}, \dots, \bar{\lambda}_{T-1}\}$.

For example, $\bar{\lambda}_B$ is the FOI level experienced continuously from the start of year B until its end.

When an FOI is constant within a year, we can directly integrate Equation (8) between the start of a year and its end. We illustrate this between the year of birth, B , and the following year; in doing so, we assume individuals are born at the start of year B :

$$\int_B^{B+1} \frac{d}{dt'} \left(X^B(t') \exp(\bar{\lambda}_B t') \right) dt' = \int_B^{B+1} \bar{\lambda}_B \exp(\bar{\lambda}_B t') dt' \quad (14)$$

which yields the following solution for the proportion seropositive one year after their birth:

$$X^B(B+1) = 1 - \exp(-\bar{\lambda}_B) \quad (15)$$

Repeating this for the next year, we have:

$$\begin{aligned} X^B(B+2) \exp(\bar{\lambda}_{B+1}(B+2)) - X^B(B+1) \exp(\bar{\lambda}_{B+1}(B+1)) \\ = \exp(\bar{\lambda}_{B+1}(B+2)) - \exp(\bar{\lambda}_{B+1}(B+1)) \end{aligned} \quad (16)$$

which can be rearranged to:

$$X^B(B+2) = 1 + (X^B(B+1) - 1) \exp(-\bar{\lambda}_{B+1}) \quad (17)$$

$$= 1 - \exp(-(\bar{\lambda}_B + \bar{\lambda}_{B+1})) \quad (18)$$

Iteratively, for an arbitrary year $T > B$, we find that:

$$X^B(T) = 1 - \exp\left(-\sum_{i=B}^{T-1} \bar{\lambda}_i\right) \quad (19)$$

Equation (19) effectively defines a survival model, where the probability of becoming infected by a given time relates to the cumulative force of infection to which an individual has been exposed. If there is strong transmission throughout an individual's life, they are more likely to have been infected and be seropositive by the time of antibody measurement.

It is also possible to parameterize Equation (19) in terms of a birth-cohort-specific age: $A^B = T - B$ (in integer years), where B is fixed and T varies:

$$X^B(A^B) = 1 - \exp\left(-\sum_{i=0}^{A^B-1} \bar{\lambda}_i^B\right) \quad (20)$$

where $\bar{\lambda}_A^B := \bar{\lambda}_{B+A}$ is the force of infection experienced by a cohort born in year B as a function of its age. When transmission varies over time, different birth cohorts can experience varied histories of transmission. This means that each age group will experience its own set of $\bar{\lambda}_A^B$ values. For example, suppose that those aged one at the start of 2024 (i.e., those born at the start of 2023) were subject to a constant force of infection $\bar{\lambda}_{2023}$ in their first year of life. This means the proportion seropositive will be:

$$X^{2023}(A^{2023} = 1) = 1 - \exp(-\bar{\lambda}_{2023}) \quad (21)$$

Those aged two at the start of 2024 will experience the same force of infection in their second year as those aged one did in their first, $\bar{\lambda}_{2023}$. But they will have also experienced a force of infection

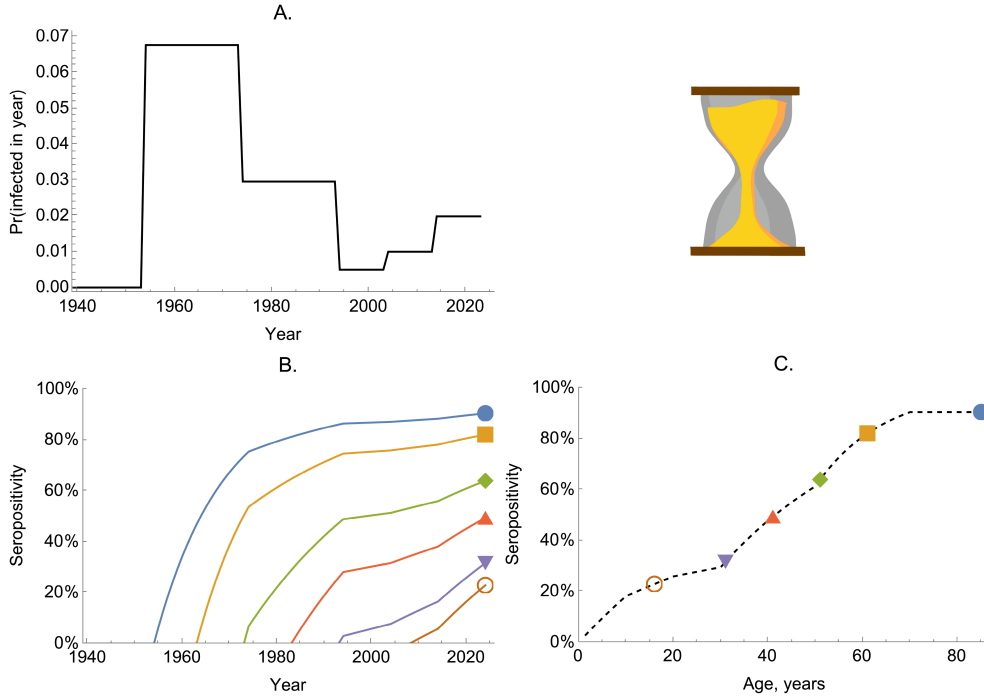


FIGURE 3 | The dynamics of seropositivity in a time-varying FOI model. Panel A shows the probability of becoming infected in a given year (given by $1 - \exp(-\bar{\lambda}_T)$) over time. Panel B shows the solution of Equation (19) for six birth cohorts (colored lines). Panel C shows the proportion seropositive by age for the population in 2024; the dashed line shows the solution for all age groups, and the colored markers correspond to the seropositivity shown in 2024 in Panel B.

λ_{2022} in their first year of life, meaning the proportion seropositive of this cohort at the start of 2024 is given by:

$$X^{2022}(A^{2022} = 2) = 1 - \exp(-(\bar{\lambda}_{2022} + \bar{\lambda}_{2023})) \quad (22)$$

We now consider the seropositivity profile in a given year T as a function of integer age A of individuals within a population: This emulates the typical data obtained from serological surveys, which can be obtained from Equation (19):

$$X_T(A_T) := 1 - \exp\left(-\sum_{i=0}^{A_T-1} \bar{\lambda}_i^{-T-A_T}\right) \quad (23)$$

where $A_T = T - B$ corresponds to the age of a cohort born at the start of year B . We use the notation A_T to denote the age of a member of the population surveyed in year T , that is, B varies across the population; this is distinguished from A_T which is the age of a particular birth cohort at time T , that is, B is fixed.

In Figure 3, we illustrate the impact of a time-varying FOI by considering a range of birth cohorts for a pathogen with varied historical transmission rates (Panel A). Panel B shows the solution given by Equation (19) for six different age cohorts: In each case, we solve the model from their birth until the present time (taken as $T = 2024$ in Equation (19)). In Panel C, we show the seropositivity of each age cohort in 2024 as calculated by Equation (23).

5.1 | Including Seroreversion

Until now, we have assumed that antibodies, as a signature of past infection, are always detectable. For some infections, however,

antibody levels wane with time below the limit of detection and *seroreversion* is said to have occurred, defined in this context as a confirmed positive serologic test later testing negative [11, 12]. We now discuss how to incorporate seroreversion into serocatalytic models by building on the time-varying model introduced above.

We modify Equation (3) to allow seropositive individuals to become seronegative:

$$\begin{aligned} \frac{dS^b(t)}{dt} &= -\lambda(t)S^b(t) + \mu X^b(t), \\ \frac{dX^b(t)}{dt} &= \lambda(t)S^b(t) - \mu X^b(t), \\ S^b(b) &= 1, X^b(b) = 0 \end{aligned} \quad (24)$$

and assume that the rate of seroreversion, $\mu > 0$, is constant. Since $S^b(t) + X^b(t) = 1$ for $t \geq b$, we can rewrite the equation for seropositive individuals:

$$\frac{dX^b(t)}{dt} = \lambda(t)(1 - X^b(t)) - \mu X^b(t) \quad (25)$$

We then separate the variables:

$$\frac{dX^b(t)}{dt} + (\mu + \lambda(t))X^b(t) = \lambda(t) \quad (26)$$

When λ is constant over time (i.e., $\lambda(t) = \bar{\lambda}$), we can solve Equation (26) for the proportion seropositive at time $t \geq b$ using the integrating factor approach:

$$X^b(t) = \frac{\bar{\lambda}}{\lambda + \mu} [1 - \exp(-(\bar{\lambda} + \mu)(t - b))] \quad (27)$$

Equation (27) shows that the proportion seropositive increases monotonically towards a plateau at $\frac{\bar{\lambda}}{\bar{\lambda} + \mu}$, which is below 1 (because $\mu > 0$).

We can alternatively assume that $\lambda(t)$ is piecewise-constant, where each calendar year has a given (constant) FOI. Considering the first year of life for a birth-cohort born at the start of year B , we can determine the proportion seropositive by using an integration factor approach:

$$\begin{aligned} & \int_B^{B+1} \frac{d}{dt'} (X^B(t') \exp((\mu + \bar{\lambda}_B)t')) dt' \\ &= \int_B^{B+1} \bar{\lambda}_B \exp((\mu + \bar{\lambda}_B)t') dt' \end{aligned} \quad (28)$$

When rearranged, this yields the following expression for the proportion seropositive at the start of year $B + 1$:

$$X^B(B + 1) = \frac{\bar{\lambda}_B}{\mu + \bar{\lambda}_B} (1 - \exp(-(\mu + \bar{\lambda}_B))) \quad (29)$$

Repeating this exercise for an arbitrary year, $T > B$, produces the following:

$$X^B(T + 1) = \frac{\bar{\lambda}_T}{\mu + \bar{\lambda}_T} + \left(X^B(T) - \frac{\bar{\lambda}_T}{\mu + \bar{\lambda}_T} \right) \exp(-(\mu + \bar{\lambda}_T)) \quad (30)$$

If $\mu = 0$, then Equation (30) can be used to produce to Equation (19).

Equation (30) is an iterative solution, with the proportion seropositive at the end of each year becoming the initial condition for the next piece. The symbolic expression for the solution that results from this iterative process is cumbersome, but it is straightforward to solve for a general $X^B(T)$, with $T \geq B$, numerically using a for loop (see Algorithm 1).

In Figure 4, we show how seropositivity profiles are modified by incorporating seroreversion. Unlike the model without seroreversion, seropositivity of birth cohorts can decrease over time (Figure 4B): this occurs when the proportion of seropositive

individuals becoming seronegative exceeds the proportion of seronegative individuals becoming infected, that is, if:

$$\mu X^b(t) > \lambda(1 - X^b(t)) \quad (31)$$

which gives a threshold condition: $\frac{\lambda}{\lambda + \mu} < X^b(t)$.

Overall, the inclusion of seroreversion in the model leads to lower seropositivity at each age (Figure 4C). In Figure 4B, we see that older cohorts tend to have higher seropositivity than younger cohorts. We ask if this is generally true for time-varying FOI models with seroreversion. The answer is *nearly*: Older cohorts must have seropositivity *at least* as great as younger cohorts. When a new cohort is born, an existing cohort must have a seropositivity of zero or higher. From then on, the two cohorts experience the same forces of infection, and differentiating Equation (30) with respect to $X^B(T)$, we find:

$$\frac{\partial X^B(T + 1)}{\partial X^B(T)} = \exp(-(\bar{\lambda}_T + \mu)) > 0 \quad (32)$$

meaning that the cohort starting an interval of fixed $\bar{\lambda}_T$ with seropositivity greater than another will, throughout the interval, have higher seropositivity. This means that the separate lines representing seropositive fractions through time shown in Figure 4B will never meet.

5.2 | Data Example: Chikungunya Virus

Chikungunya virus is widespread in the tropics, where it causes recurrent outbreaks of chikungunya fever. The virus is transmitted by *Aedes aegypti* and *Aedes albopictus* mosquitoes [13]. Chikungunya fever is characterized by severe arthralgia and myalgia that can persist for years and have considerable detrimental effects on health, quality of life, and economic productivity [13]. Clinical presentation of chikungunya fever is often non-discernible from that of other arboviruses, and it can also be confused with other febrile illnesses [13]. This means that a chikungunya epidemic is likely to go unnoticed or missed, and therefore, serosurveys are instrumental in understanding the disease burden or virus epidemiology.

ALGORITHM 1 | R function for solving the time-varying model.

```
solve_time_varying_model <- function(T, B,
  lambdas, mu) {
  # initial condition for X
  X = c(0)
  age = T - B

  for (i in 1:age) {
    T_i = B + i
    # lambdas[T_i] corresponds to the FOI value in year T_i
    X[i + 1] = (lambdas[T_i] / (mu + lambdas[T_i])) +
      (X[i] - (lambdas[T_i] / (mu + lambdas[T_i]))) * exp(- (mu + lambdas[T_i]))
  }
  return(X)
}
```

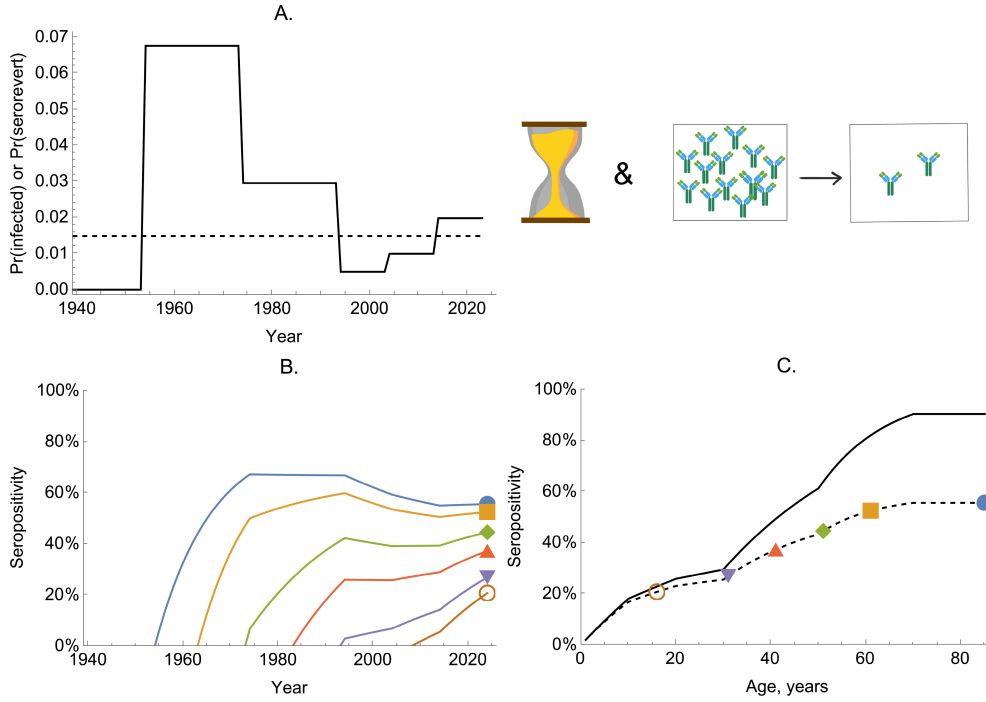


FIGURE 4 | The dynamics of seropositivity in a model with time-varying FOI and seroreversion. Panel A shows the probability of becoming infected per year (given by $1 - \exp(-\bar{\lambda}_T)$) over time (black solid line) and the fixed probability of seroreversion per year (dashed line). Panel B shows the proportion seropositive for six birth cohorts (colored lines). Panel C shows the proportion seropositive by age for the population in 2024; the dashed line shows the proportion seropositive for all age groups for the model including seroreversion and the colored markers correspond to the seropositivity shown in 2024 in Panel B; the solid line shows the model solution if the rate of seroreversion were zero.

Here, we analyze CHIKV serological data derived from cross-sectional surveys undertaken in 2015 to estimate historical transmission patterns in Burkina Faso and Gabon [14]. This dataset includes serological information for individuals between 1 and 55 years old (see Figure 5A). Chikungunya is an epidemic disease, and to account for this, we fit a time-varying serocatalytic model (in a Bayesian framework, see Section A), though in this case the model does not allow for seroreversion. Our model provided a reasonable fit to the serological data (Figure 5A). Our reconstructed series for the annual probability of infection (Figure 5B) indicated a likely reduction in transmission more recently, which is similar to that determined in previous work [14].

6 | Age-Dependent Force of Infection

We now explore the case where FOI varies by age, but the profile of infection remains constant over time. That is, the FOI experienced by a 5-year-old in 1939 is the same as that experienced by someone of the same age in 2024 but distinct from those aged 10.

6.1 | Without Seroreversion

We first suppose there is no seroreversion, meaning the dynamics of seropositivity are described by a slight modification of Equation (4):

$$\frac{dX(a)}{da} = \lambda(a)(1 - X(a)), \quad X(a=0) = 0 \quad (33)$$

where $a = t - b$ is the age at time t of serosurvey for the cohort born at time b . We note, however, that Equation (33) uses a different notation to Equation (4), where we no longer use X^b ; rather, we use X , because, if we know the age of individuals, when they were born does not influence the trajectories of their seropositivity.

When FOIs are independent of age, the solution becomes:

$$X(a) = 1 - \exp(-a\bar{\lambda}) \quad (34)$$

which is the same as Equation (12), because when FOIs do not vary (through time or by age), the time-dependent model and the age-dependent FOI models become equivalent.

Generally, when we assume the FOI varies by year of age (but is constant within each year), the FOI experienced by a cohort of integer age A is given by: $\{\bar{\lambda}_0, \bar{\lambda}_1, \dots, \bar{\lambda}_{A-1}\}$. For example, a 3-year-old individual would have experienced three distinct FOIs, each corresponding to one year of their life. In this case, the solution becomes:

$$X(A) = 1 - \exp\left(-\sum_{i=0}^{A-1} \bar{\lambda}_i\right) \quad (35)$$

This means that the saturating level of seropositivity is given by: $X(\infty) = 1 - \exp\left(-\sum_{i=0}^{\infty} \bar{\lambda}_i\right) \leq 1$.

In Figure 6, we consider a pathogen with transmission patterns similar to those of sexually transmitted infections, which have a strong age-related exposure profile. Here, we model the

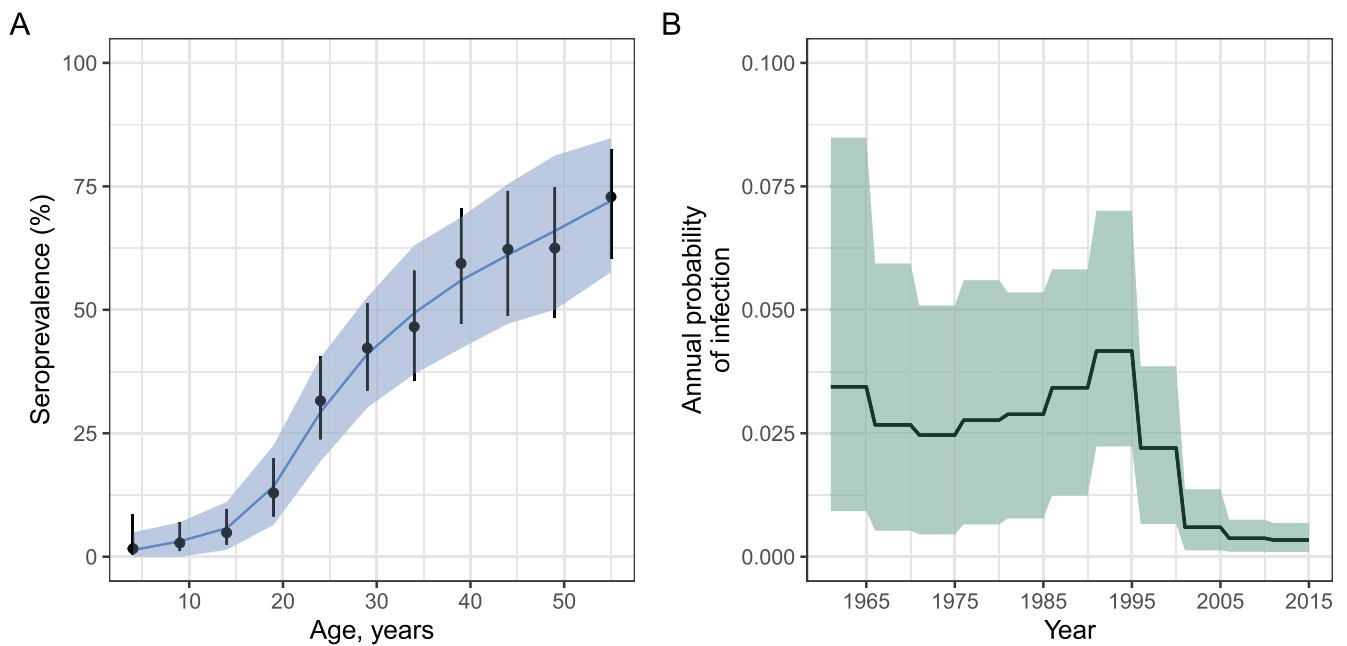


FIGURE 5 | Explaining chikungunya serological data in Burkina Faso and Gabon using a time-varying FOI model. Panel A shows the observed and fitted seroprevalence by age from surveys undertaken in 2015. Points and whiskers represent the observed proportions with 95% confidence intervals. The solid blue line indicates the mean of the posterior samples. Panel B shows the posterior mean annual probability of infection estimates given by $(1 - \exp(-\bar{\lambda}_T))$. In both panels, the shading indicates the 95% credible intervals, representing the 2.5th and 97.5th percentiles of the posterior distributions. We assumed that the FOI was piecewise-constant with pieces of width 5 years.

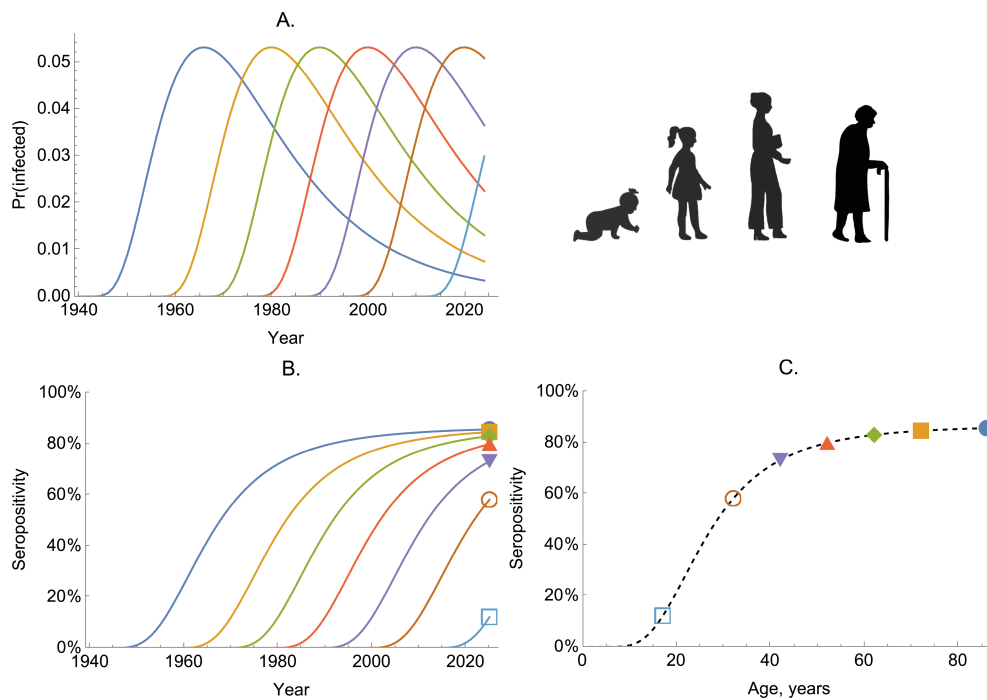


FIGURE 6 | The dynamics of seropositivity with age-dependent infection risk—serological dynamics of a sexually transmitted infection. Panel A shows the probability of becoming infected per year (given by $1 - \exp(-\bar{\lambda}_A)$) for seven birth cohorts (colored lines). Panel B shows the proportion seropositive for the same birth cohorts (colored lines). Panel C shows the proportion seropositive by age for the population in 2024; the dashed line shows the proportion seropositive for all age groups; the colored markers correspond to the same 2024 values as shown at the right edge of Panel B.

serological dynamics of such a disease where the FOIs are only dependent on age and not time, that is, we assume that transmission has been stable through time. Specifically, we assume that FOI is highest in those aged in their early 20s and that the age

profile of infection risk remains constant over time (Figure 6A). The trajectories of seropositivity follow the same profile, regardless of when a cohort was born (Figure 6B), and older cohorts have higher seropositivity (Figure 6C).

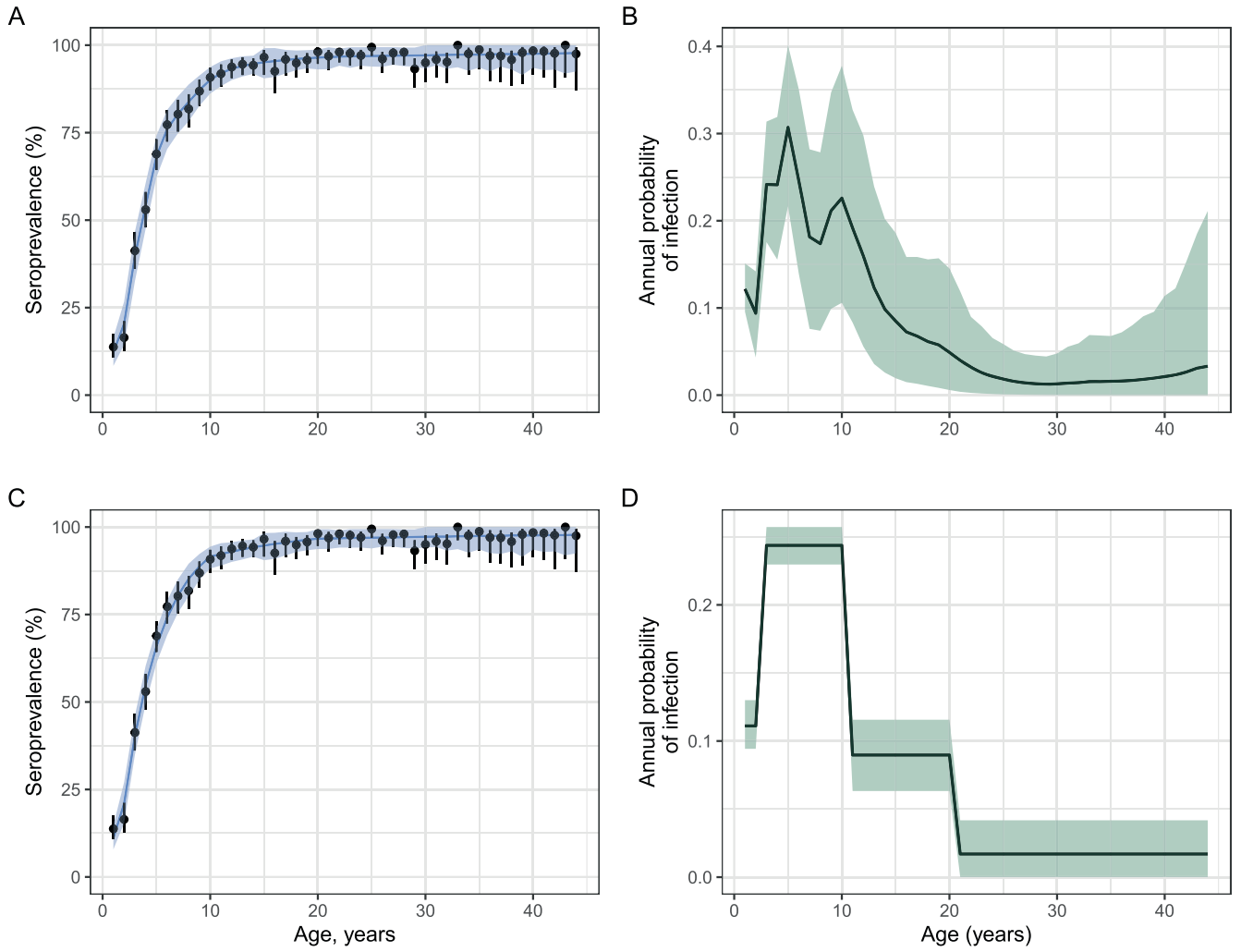


FIGURE 7 | Age-varying FOI model fits for mumps virus. Panels A and C show the observed and model-fitted seroprevalence by age. Points and whiskers represent the observed proportions with 95% confidence intervals. The blue lines indicate the mean of posterior samples, and the shading shows the 95% credible intervals representing the 2.5th and 97.5th percentiles of the posterior distributions. Panels B and D show the estimated age-specific annual probability of infection, calculated as $1 - \exp(-\bar{\lambda}_A)$, from a model with one-year piecewise-constant FOIs (Panel B) and a model assuming four FOI pieces (Panel D). For Panels B and D, we estimated four FOI values across the following age bands: (1,2), (3–8), (9,10), (11–44).

6.2 | With Seroreversion

We now consider an age-dependent FOI model with seroreversion using a slight modification to Equation (33):

$$\frac{dX(a)}{da} = \lambda(a)(1 - X(a)) - \mu X(a), \quad X(a=0) = 0 \quad (36)$$

We solve this equation assuming FOI is piecewise-constant with yearly piece-widths. To do so, we denote integer-valued age by an uppercase A :

$$\frac{dX(a')}{da'} = \bar{\lambda}_A(1 - X(a')) - \mu X(a'), \quad a' \in [A, A+1] \quad (37)$$

An integration factor approach can be used to yield:

$$\begin{aligned} \int_A^{A+1} \frac{d}{dt'} (X(a') \exp((\mu + \bar{\lambda}_A)a')) da' \\ = \int_A^{A+1} \bar{\lambda}_A \exp((\mu + \bar{\lambda}_A)a') da' \end{aligned} \quad (38)$$

This has the same solution as given by Equation (30) (replacing $T \rightarrow A$):

$$X(A+1) = \frac{\bar{\lambda}_A}{\bar{\lambda}_A + \mu} + \left(X(A) - \frac{\bar{\lambda}_A}{\bar{\lambda}_A + \mu} \right) \exp(-(\bar{\lambda}_A + \mu)) \quad (39)$$

where $X(A+1) > X(A)$ if $X(A) < \frac{\bar{\lambda}_A}{\bar{\lambda}_A + \mu}$.

6.3 | Data Example: Mumps Virus

Mumps is a common childhood infection caused by the mumps virus. The hallmark of infection is swelling of the parotid gland, and aseptic meningitis and encephalitis are common complications. Other complications include deafness and pancreatitis [15]. We analyzed cross-sectional seroprevalence data collected in 1986–1987 in the UK for individuals between 1 and 44 years of age (Figure 7A,C) [16]. We assumed that the FOI was age-varying as in the reporting study [16].

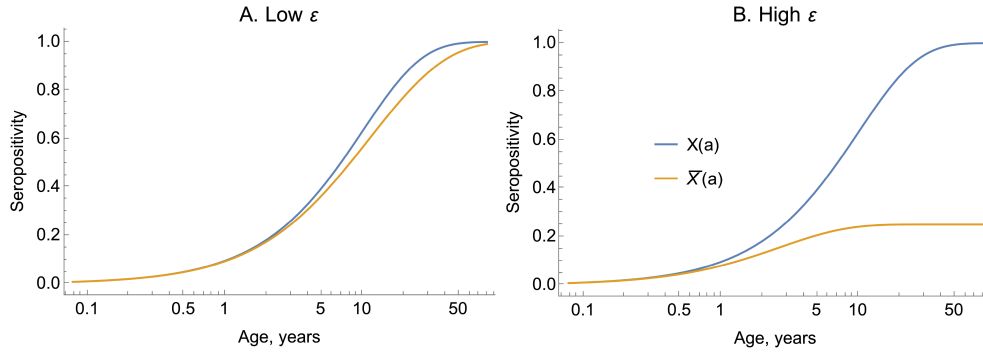


FIGURE 8 | An elevated death rate subsequent to infection can dilute the pool of seropositive individuals. In both panels, we plot seropositivity estimated using Equation (42) (orange lines), which accounts for an elevated death rate post-infection. In Panel A, we assume a low death rate, $\epsilon = 0.05$; in Panel B, we assume a higher death rate, $\epsilon = 0.4$. In both panels, $\lambda = 0.1$. The blue lines show the approximate solution given by neglecting to account for the elevated death rate (Equation (13)).

We considered two scenarios: (i) piecewise-constant FOIs with one-year age widths, that is, a distinct FOI for each age class, (ii) four FOI change points corresponding to ages 2, 8, 10, and 24 years. The first scenario is a more flexible model assumption but prone to wide uncertainty around each λ_A since there are relatively few data points to inform each FOI. The second scenario characterizes transmission intensity across various age-bins—for example, it allows different FOIs in infancy, during pre-school years, and during adolescence or teenage years. This scenario may result in lower variance estimates but is at a higher risk of bias.

In both scenarios, the model fits were similar and reasonable (Figure 7A,C). In both models, the inferred FOI peaked in younger children and was low for older age groups. The estimated age-specific annual probability of infection was highest in the age range 2–10 years, with a peak at 0.36 (95% CrI: 0.24–0.52) in scenario (i) (Figure 7B) and 0.25 (95% CrI: 0.26–0.30) in scenario (ii) (Figure 7D). The distinct model fits illustrate the impact of assumptions during inference as shown by higher variance in FOIs in model (i) (Figure 7C).

7 | Elevated Death Rate Due to Infection

We now suppose individuals exposed to a pathogen have an elevated risk of death. We show that exposure-related mortality biases the observed serological profile and impacts how data should be interpreted. We illustrate this effect by including an additional “deceased (D)” compartment to the age-dependent FOI model:

$$\begin{aligned} \frac{dS(a)}{da} &= -\lambda(a)S(a), \\ \frac{dX(a)}{da} &= \lambda(a)S(a) - \epsilon X(a), \\ \frac{dD(a)}{da} &= \epsilon X(a) \end{aligned} \quad (40)$$

where $\epsilon > 0$ denotes the rate at which previously infected individuals die, and we assume $S(0) = 1$, $X(0) = 0$, $D(0) = 0$ as initial conditions for each birth cohort. If we consider a special case of Equation (40) for when FOI is time-constant, we can exactly solve the system to give:

$$\begin{aligned} S(a) &= \exp(-\bar{\lambda}a), \\ X(a) &= \frac{\bar{\lambda}}{\bar{\lambda} - \epsilon} (\exp(-\epsilon a) - \exp(-\bar{\lambda}a)), \\ D(a) &= \frac{\bar{\lambda}(1 - \exp(-\epsilon a)) - \epsilon(1 - \exp(-\bar{\lambda}a))}{\bar{\lambda} - \epsilon} \end{aligned} \quad (41)$$

Generally, serological data are collected from survivors. This means that the proportion of those seropositive out of those alive is given by the ratio: $\bar{X}(a) := X(a)/(S(a) + X(a))$. Using Equation (41), we can derive a simplified expression for this quantity as follows:

$$\bar{X}(a) = 1 - \frac{\bar{\lambda} - \epsilon}{\bar{\lambda} \exp((\bar{\lambda} - \epsilon)a) - \epsilon} \quad (42)$$

If $\bar{\lambda} > \epsilon$, as $a \rightarrow \infty$, then $\bar{X}(a) \rightarrow 1$; if $\bar{\lambda} < \epsilon$, $\bar{X}(a) \rightarrow \bar{\lambda}/\epsilon < 1$. Figure 8 illustrates the analytical solution in Equation (42) in two scenarios: (i) small ϵ , (ii) versus a large ϵ .

In Figure 8A, we show that when ϵ is small then the exact solution of Equation (42) is nearly equal to the approximate solution which neglects deaths due to infection, $\bar{X}(a) = 1 - \exp(-\bar{\lambda}a)$. Conversely, when ϵ is large, Equation (42) results in lower $\bar{X}(a)$ values (orange line, Figure 8B).

7.1 | Subpopulations With Differing Mortality Risks

An issue with the model described in Section 7 is that it supposes all those who become infected will die, eventually, due to their infection. However, this is not generally true—usually only a subset of those infected will die due to their infection. We can represent the heterogeneity in the risk of severe infection leading to death by separating the seropositive populations into two groups: one, X_m , comprising individuals experiencing a *mild* infection that does not lead to an elevated death rate; and another, X_s , with individuals who experience severe infection and eventually succumb to the infection. The modified model system is then given by:

```

solve_subpopulations_mortality_model <- function(age, lambdas, epsilon, rho){
  S = c(1)
  Xm = c(0)
  Xs = c(0)

  for (i in 2:(age + 1)){
    S[i] = S[i - 1] * exp(-lambdas[i-1])
    Xm[i] = Xm[i - 1] + S[i - 1] * (1 - rho) * (1 - exp(-lambdas[i-1]))
    Xs[i] = ((S[i - 1] * rho * lambdas[i-1]) * (exp(-epsilon) - exp(-lambdas[i-1])) +
             Xs[i - 1] * (lambdas[i-1] - epsilon) * exp(-epsilon)) / (lambdas[i-1] - epsilon)
  }

  return(list(S = S, Xm = Xm, Xs = Xs))
}

```

$$\begin{aligned}
\frac{dS(a)}{da} &= -\lambda(a)S(a), \\
\frac{dX_m(a)}{da} &= (1 - \rho)\lambda(a)S(a), \\
\frac{dX_s(a)}{da} &= \rho\lambda(a)S(a) - \epsilon X_s(a), \\
\frac{dD(a)}{da} &= \epsilon X_s(a)
\end{aligned} \tag{43}$$

where $0 \leq \rho \leq 1$ represents the proportion of infections which are severe (leading to death); alternatively, known as the *infection fatality ratio*.

In this system, the proportion of those living who are seropositive, assuming λ is constant, is given by:

$$\begin{aligned}
\bar{X}(a) &:= \frac{X_m(a) + X_s(a)}{X_m(a) + X_s(a) + S(a)}, \\
&= 1 - \frac{\bar{\lambda} - \epsilon}{(1 - \rho)(\bar{\lambda} - \epsilon) \exp(\bar{\lambda}a) + \rho(\bar{\lambda} \exp((\bar{\lambda} - \epsilon)a) - \epsilon)} \tag{44}
\end{aligned}$$

If $\rho = 1$, Equation (44) reduces to Equation (42). If $\rho = 0$, the system becomes equivalent to age-dependent FOI models without death.

Assuming λ is piecewise-constant, with one-year pieces, we can determine the proportions susceptible and with mild (m) or severe (s) infection histories at the end of a given piece by:

$$\begin{aligned}
S(a+1) &= S(a) \exp(-\bar{\lambda}_a), \\
X_m(a+1) &= X_m(a) + S(a)(1 - \rho) \left(1 - \exp(-\bar{\lambda}_a)\right), \\
X_s(a+1) &= \frac{S(a)\rho\bar{\lambda}_a(\exp(-\epsilon) - \exp(-\bar{\lambda}_a)) + X_s(a)(\bar{\lambda}_a - \epsilon) \exp(-\epsilon)}{\bar{\lambda}_a - \epsilon} \tag{45}
\end{aligned}$$

Equation (45) can then be used numerically to determine the proportion seropositive of those surviving (i.e., those sampled in a serosurvey). Similarly, Equation (45) is an iterative solution, with the proportion seropositive at the end of each year of life

becoming the initial condition for the next piece as shown in Algorithm 2.

7.2 | Data Example: Ebola Virus

The West African Ebola virus disease outbreak of 2014–16 spread from Guinea to Liberia, Sierra Leone, Mali, Senegal, and Nigeria, and in 28 months, the outbreak had resulted in 28 652 cases and 11 325 deaths [17]. We fit the model described in Section 7.1 to antibody prevalence data from Ebola survivors in Sierra Leone after the 2014–2016 outbreak [18]. The Sierra Leone study was conducted between 2016 and 2018 and measured IgG to Ebola virus glycoprotein in 1282 individuals of median age 16 years, with an inter-quartile range of 7–25 years, of whom 107 (8.4%) were seropositive. In this example, we assumed a stepped time-varying FOI, with $\lambda = 0$ before the outbreak in 2014 and an estimated λ from 2014 onwards.

For one set of model fits, we set the infection fatality ratio, $\rho = 0.89$ [19], and in another $\rho = 0$. Both models, with and without the assumption of death due to infection, fitted the observed data reasonably (Figure 9A) but estimated substantially different FOIs. The model that incorporated elevated death due to infection estimated an annual probability of infection around six times higher than the model neglecting this (Figure 9B). This example illustrates the importance of accounting for infection-induced mortality in an appreciable fraction of cases—failing to do so can severely underestimate the level of transmission.

8 | Time- and Age-Dependent Force of Infection

Exposure and susceptibility can vary both by age and over calendar time. That is, $\lambda = \lambda(a, t)$. Here, we consider a simplified case where the FOI varies as a product of time- and age-specific patterns, $\lambda(a, t) = u(a)v(t)$, implying persistent age-related patterns that fluctuate temporally. A sexually transmitted infection might have FOI following such patterns, where sexual contact rates would typically peak in the early twenties, but the transmission rate may be affected by public health policies or treatment and

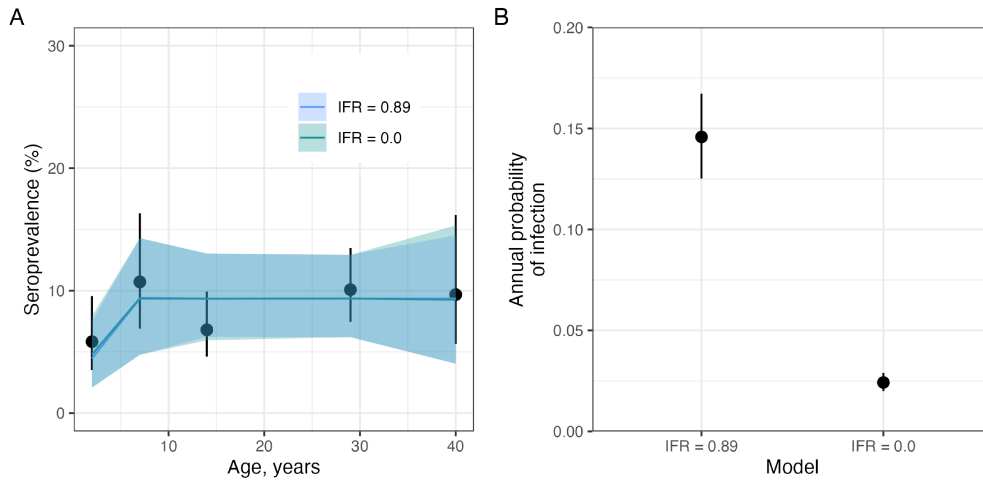


FIGURE 9 | Accounting for an elevated death rate due to infection results in a higher inferred FOI for Ebola virus disease in Sierra Leone during the West African Ebola outbreak in 2014–16. Panel A shows the observed (black circles) seroprevalence with 95% confidence intervals and model-predicted seroprevalence by age (blue and green lines). The predicted seroprevalence (blue and green lines) is derived from two models with and without the infection fatality ratio (IFR). Panel B shows the annual probability of infection for the two IFR models.

preventative measures, lowering the overall transmission rates across all ages [20].

We first consider models without seroreversion, where the proportion seropositive evolves according to:

$$\begin{aligned} \frac{dX^b(t)}{dt} &= u(t-b)v(t)(1-X^b(t)), \\ X^b(t=b) &= 0 \end{aligned} \quad (46)$$

As previously, we assume that $u(\cdot)$ and $v(\cdot)$ are both piecewise-constant with pieces of width one year; we denote their levels within integer year T for a cohort born at the start of year B by $\bar{u}(T-B)$ and $\bar{v}(T)$. This means that within a year, their product, λ is constant, and the proportion seropositive at the end of the year is given by:

$$X^B(T+1) = 1 + (X^B(T) - 1) \exp(-\bar{u}(T-B)\bar{v}(T)) \quad (47)$$

We can then solve for the proportion seropositive at the start of an arbitrary year $T > B$ for a birth-cohort born at the start of year B :

$$X^B(T) = 1 - \exp\left(-\sum_{i=0}^{T-B-1} \bar{u}(i)\bar{v}(B+i)\right) \quad (48)$$

To exemplify the behavior of this model, we consider a disease with an age-dependent FOI structure similar to that of a sexually transmitted infection: With $u(\cdot)$ peaking in the early 20s. At the same time, we allow a rapid increase in the level of transmission, starting in the year 1980 (Figure 10A), for example, representing a widespread increase in the number of infected individuals. This upward shift in transmission means that the individuals with peak sexual activity during the period of elevated transmission experience a disproportionately higher risk of infection, and the seropositive proportion of these cohorts can exceed those of older individuals (Figure 10B,C).

8.1 | Including Seroreversion

If we allow seroreversion, the seropositivity at the start of year T is given by (through analogy to Equation (30)):

$$\begin{aligned} X^B(T+1) &= \frac{\bar{u}(T-B)\bar{v}(T)}{\bar{u}(T-B)\bar{v}(T) + \mu} \\ &+ \left(X^B(T) - \frac{\bar{u}(T-B)\bar{v}(T)}{\bar{u}(T-B)\bar{v}(T) + \mu} \right) \\ &\times \exp(-(\bar{u}(T-B)\bar{v}(T) + \mu)) \end{aligned} \quad (49)$$

Like for the time- and age-only models, when $u(\cdot)$ and $v(\cdot)$ vary by age and year, there is no simple analytical expression for the seropositivity. Nonetheless, Equation (49) can be used to produce an iterative solution if we assume $u(\cdot)$ and $v(\cdot)$ are piecewise-constant.

8.2 | Data Example: HIV

The force of infection for HIV likely varies over time and by age. Additionally, the majority of infected individuals, bar a few rare cases, remain infected throughout their lives. As a result, we can model HIV infection prevalence in the same way we model seropositivity for other pathogens, using models without seroreversion. In the absence of treatment, HIV infection eventually leads to death, and so it is crucial to account for infection-related mortality when analyzing HIV prevalence data. Additionally, it is important to account for the delay from infection to death. Progression to AIDS and death typically follows a long asymptomatic period—estimated to be around 10 years [21].

To model HIV serodynamics, we use the following system:

$$\begin{aligned} \frac{dS^b(t)}{dt} &= -u(t-b)v(t)S^b(t), \\ \frac{dX_1^b(t)}{dt} &= u(t-b)v(t)S^b(t) - X_1^b(t), \end{aligned}$$

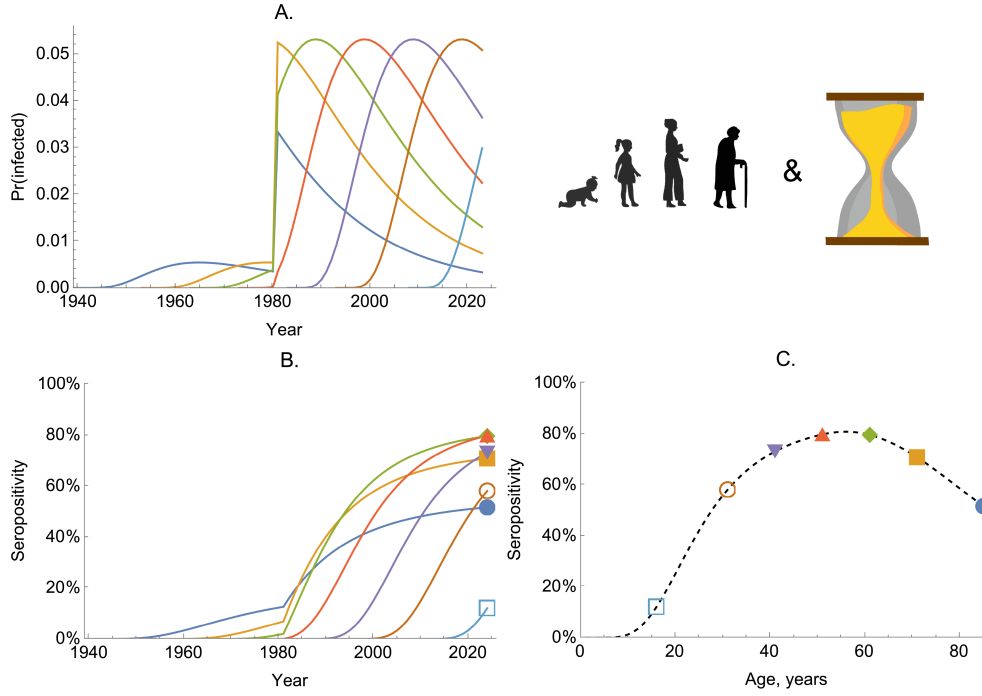


FIGURE 10 | The dynamics of seropositivity with age- and time-dependent infection risk. Panel A shows the probability of becoming infected per year (given by $1 - \exp(-\bar{\lambda})$) for seven birth cohorts (colored lines). Panel B shows the proportion seropositive for the same birth cohorts (colored lines). Panel C shows the proportion seropositive by age in a single cross-section, taken in 2024. The dashed line shows the proportion seropositive across all age groups, while the colored markers correspond to the same values shown in 2024 in Panel B.

$$\begin{aligned} \frac{dX_i^b(t)}{dt} &= X_{i-1}^b(t) - X_i^b(t), \quad \text{for } i = 2, 3, \dots, 10, \\ \frac{dD^b(t)}{dt} &= X_{10}^b(t) \end{aligned} \quad (50)$$

where $S^b(b) = 1$ and all other compartments have initial conditions set to zero. In Equation (50), we have made use of the linear-chain trick (see, e.g., [10]) by splitting the HIV infection-positive individuals into 10 compartments: X_1, \dots, X_{10} , and death occurs only after the last compartment is reached. Because the rate parameters in front of each X_i^b term are 1 in Equation (50), this results in a typical duration to death due to infection following a gamma(10, 1) distribution, which has a mean of 10 years (see²).

As before, we assume that $u(\cdot)$ and $v(\cdot)$ are piecewise-constant with one-year pieces. Since the system is linear, we can write it as a vector differential equation by defining $\mathbf{Y}^b(t) := [S^b(t), X_1^b(t), \dots, X_{10}^b(t), D^b(t)]'$ as a vector of states and $\bar{\mathbf{A}}_{T,B}$ as a matrix of constants that are fixed for a given calendar year T assuming individuals were born at the start of year B :

$$\bar{\mathbf{A}}_{T,B} := \begin{bmatrix} -\bar{u}(T-B)\bar{v}(T) & 0 & \dots & & & \\ \bar{u}(T-B)\bar{v}(T) & -1 & 0 & \dots & & \\ 0 & 1 & -1 & 0 & \dots & \\ 0 & 0 & 1 & -1 & 0 & \dots \\ \vdots & \vdots & \ddots & \ddots & \ddots & \ddots \\ 0 & \dots & 0 & 1 & 0 & \end{bmatrix} \quad (51)$$

Using this matrix, we can write down the system of equations as a single vector equation:

$$\frac{d\mathbf{Y}^B(t)}{dt} = \bar{\mathbf{A}}_{T,B} \mathbf{Y}^B(t), \quad \text{for } t \in [T, T+1] \quad (52)$$

We can then analytically solve Equation (52) to yield the seroprevalence at the start of year $T+1$ of those born at the start of year B (note that T and B are integers):

$$\mathbf{Y}^B(T+1) = \exp(\bar{\mathbf{A}}_{T,B}) \mathbf{Y}^B(T) \quad (53)$$

where, here, $\exp(\cdot)$ denotes the matrix exponential. The matrix form means that a general solution for the system can be written down as follows:

$$\mathbf{Y}^B(T) = \exp\left(\sum_{T'=B}^{T-1} \bar{\mathbf{A}}_{T',B}\right) \mathbf{Y}^B(B) \quad (54)$$

Since all serocatalytic models described here are linear ordinary differential equation systems, we can solve such models this way; however, this abstract form generally does not permit such mathematical analysis as the more bespoke methods presented up to this point.

We fit the model structure described by Equation (50) to age-specific HIV prevalence data from a HIV survey conducted in rural Kwa Zulu Natal, South Africa, in 2003 [22] (points and uncertainty bars in Figure 11A). The survey was conducted just before the large-scale roll-out of HIV care and antiretroviral treatment (ART) in 2004 and before the substantial reduction of HIV-associated mortality and increased life expectancy [22]. The data include HIV-seroprevalence profiles of women aged 15–49 years. Older age groups did not commonly participate in the

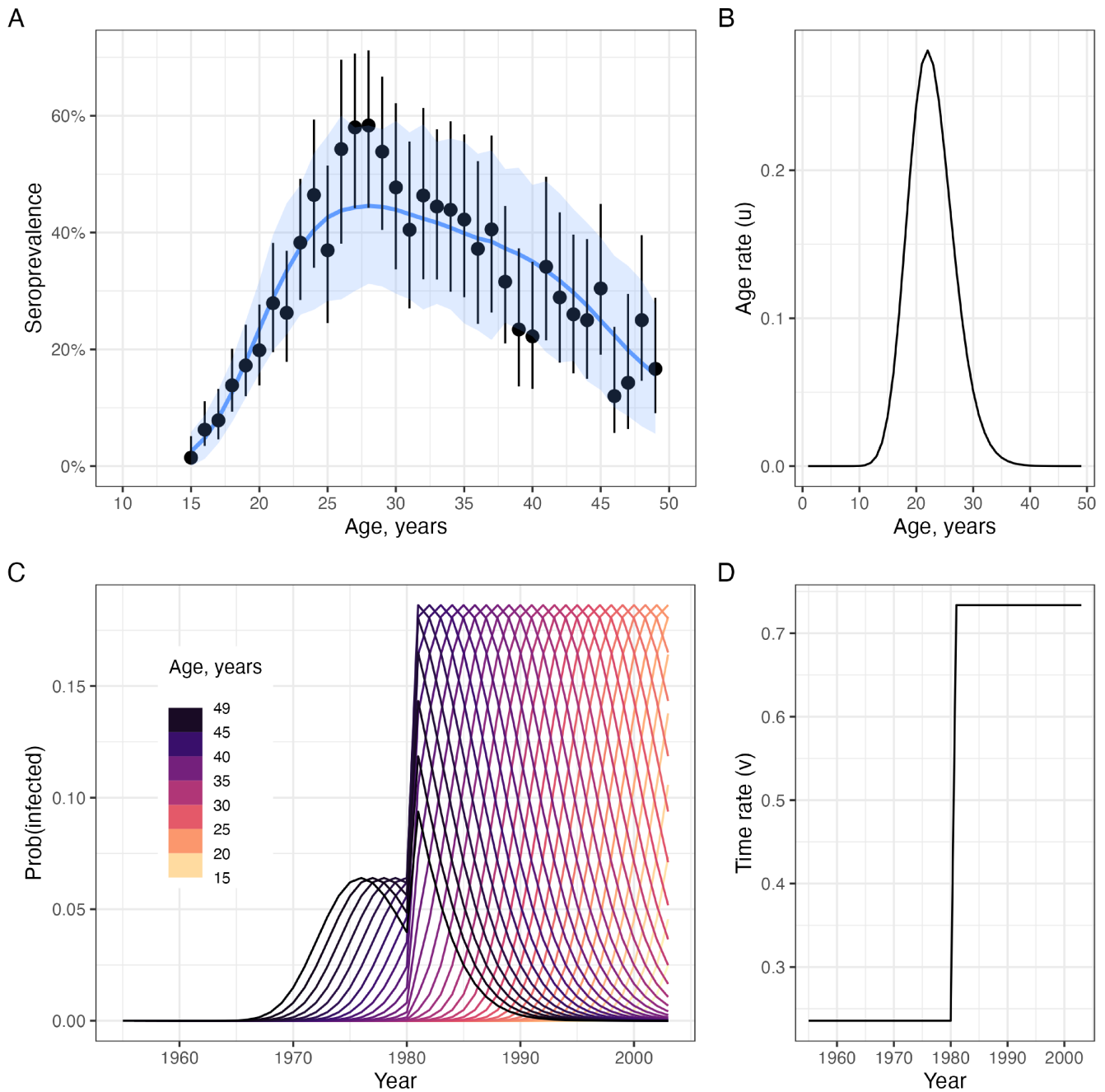


FIGURE 11 | Modeling HIV prevalence in South Africa before widespread treatment. Panel A shows the observed and model-fitted seroprevalence by age. Points and whiskers represent the observed prevalence with exact 95% confidence intervals. The blue line and shading indicate the mean of posterior samples, with 95% credible intervals representing the 2.5th and 97.5th percentiles of the model's posterior distributions. Panel B shows the estimated age patterns of infection $u(\cdot)$. Panel C shows the annual probability of becoming infected for each cohort. The probability of infection is calculated as $1 - \exp(-u(\cdot)v(\cdot))$. Panel D shows the estimated time patterns of infection $v(\cdot)$.

HIV surveillance and were only routinely included from 2007 onwards [22].

The modeled seroprevalence estimates are in broad agreement and overlap with the observed data (Figure 11A). The age of peak infection risk was estimated to lie between 20 and 30 years (Figure 11B); this is consistent with patterns of sexual behavior and partnership observed in longitudinal population-based

surveys [23]. We simplified the assumptions about changes in disease transmission over time by allowing one transmission rate multiplier prior to the year 1980 ($v(T \leq 1980)$) and another after 1980 ($v(T > 1980)$). Our model estimates indicated a roughly 3-fold increase in transmission levels between these two periods (Figure 11C,D), from $v(T \leq 1980) = 0.23$ to $v(T > 1980) = 0.73$, resulting in a commensurate increase in the probability of infection by age 40.

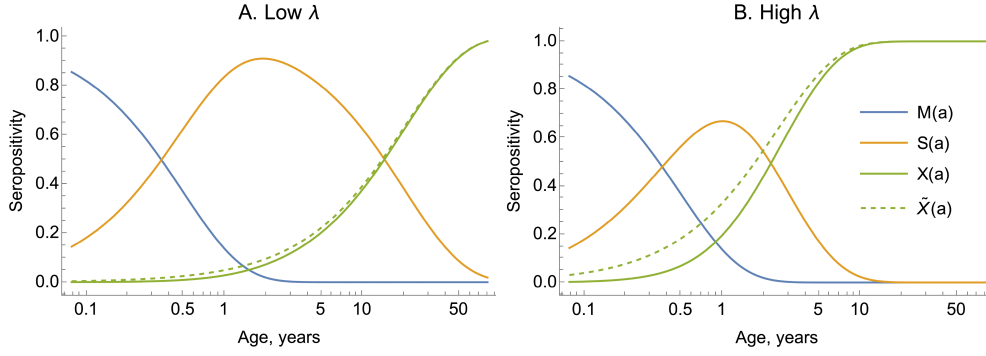


FIGURE 12 | Neglecting maternal antibody dynamics can bias the proportion seropositive. In both panels, we plot Equation (56) (solid lines), where we assume $\gamma = 2$ per year; in Panel A, we assume $\lambda = 0.05$; in Panel B, we assume $\lambda = 0.4$. The dashed lines show the approximate solution given by neglecting the effect of maternal antibodies (i.e., Equation (57)). Note, the horizontal axis is on the log scale, which cannot start at zero, and the reason the maternal compartment does not begin at $M(0) = 1$.

9 | Maternal Antibodies

For many pathogens, the offspring inherit maternal antibodies, which are typically short-lived, and yet failing to account for their influence in serological measurements can bias FOI estimates, as we now illustrate. We consider a system where FOIs are age-dependent only and with no seroreversion, and we demonstrate how the impact of maternal antibodies can be modeled. We assume that only after maternal antibody decay can individuals become naturally infected, and that the decay rate is age-independent. This results in a modified system of equations:

$$\begin{aligned} \frac{dM(a)}{da} &= -\gamma M(a), \\ \frac{dS(a)}{da} &= \gamma M(a) - \lambda(a)S(a), \\ \frac{dX(a)}{da} &= \lambda(a)S(a) \end{aligned} \quad (55)$$

where $M(a)$ is the proportion of the population with maternal antibodies of age a .

For initial conditions, we assume that $M(0) = 1$, $S(0) = 0$, $X(0) = 0$. Assuming that all births have maternal antibodies for a particular pathogen may be crude, and we discuss this in Section 9.2.

We first consider an idealized scenario when the FOI is constant. Then the above system has the following solution:

$$\begin{aligned} M(a) &= \exp(-\gamma a), \\ S(a) &= \frac{\gamma}{\gamma - \bar{\lambda}} (\exp(-\bar{\lambda} a) - \exp(-\gamma a)), \\ X(a) &= \frac{\gamma(1 - \exp(-\bar{\lambda} a)) - \lambda(1 - \exp(-\gamma a))}{\gamma - \bar{\lambda}} \end{aligned} \quad (56)$$

The rate of removal of maternal antibodies is typically high and will initially dominate the dynamics. The FOI will generally be much lower than this rate (i.e., $\bar{\lambda} \ll \gamma$), meaning $\exp(-\gamma a) \approx 0$ for older age cohorts, and the proportion seropositive will follow:

$$X(a) \approx 1 - \exp(-\bar{\lambda} a) \quad (57)$$

which reproduces Equation (34).

The typical duration of maternal antibodies is estimated to be around 6 months, meaning $\gamma = 2$ per year. In Figure 12, we consider two scenarios corresponding to a $\bar{\lambda} \ll \gamma$ (Panel A) and $\bar{\lambda} \sim \gamma$ (Panel B), where we assume $\gamma = 2$. We also plot the approximate solution given by Equation (57), which neglects the influence of maternal antibodies. This shows that when $\bar{\lambda} \ll \gamma$, it makes little difference to the modeled seropositive proportion to neglect the effect of maternal antibodies, but when $\bar{\lambda} \sim \gamma$, this induces bias in this proportion.

We typically cannot differentiate the actual proportion with maternal immunity from those with immunity due to infection during their lifetimes, which means that we can consider the overall proportion seropositive as:

$$M(a) + X(a) = 1 - S(a) = 1 - \frac{\gamma}{\gamma - \bar{\lambda}} (\exp(-\bar{\lambda} a) - \exp(-\gamma a)) \quad (58)$$

9.1 | Data Example: Enterovirus D68

Enterovirus D68 (EV-D68) infection leads to severe acute respiratory distress in children below 5 years of age with clinical symptoms of hypoxia and wheezing associated with a significant increase in pediatric hospitalizations [24]. A subset of children develop central nervous system complications, and in serious cases, respiratory failure [24]. We used EV-D68 serological data to illustrate the effect of neglecting maternal antibodies in the estimation of seroprevalence and the FOI. The data consisted of individuals aged from newborns to 80 years old, sampled in the UK in 2006 [25]. Although recent studies have suggested there may have been increased EV-D68 transmission over time, particularly after the EV-D68 emergence in 2014 [26, 27], there had not been a major EV-D68 outbreak reported in the UK at or before 2006. So, we assumed age- and time-independent risk of infection for the population sampled.

We compared two models: (i) an age- and time-independent model that considers the presence of maternal antibodies—as in Equation (56), (ii) and one that does not—Equation (57). The model in Equation (56) fitted the observed seroprevalence data well for lower age groups (ages 1–5 years) compared to the model defined in Equation (57) (Figure 13). We estimated (i) $\gamma = 1.6$ per

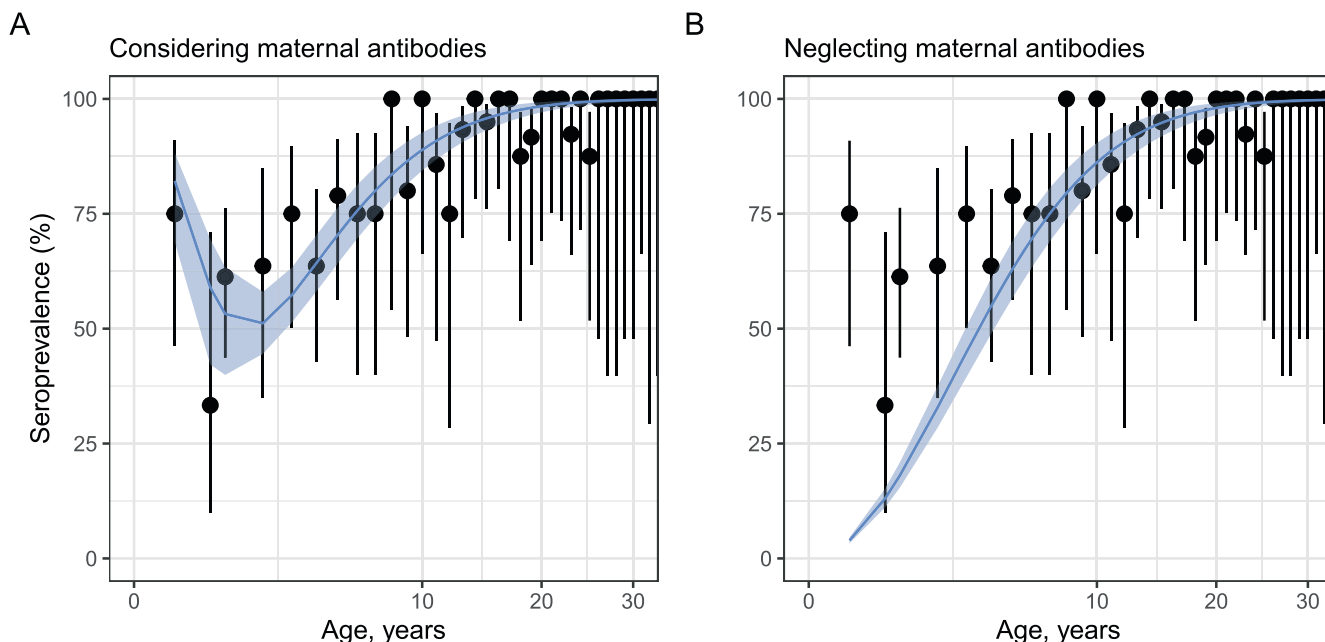


FIGURE 13 | Accounting for maternal antibody dynamics affects the inferred seroprevalence for enterovirus D68. Both panels show the observed (black points, with 2.5th and 97.5th percentiles shown as vertical error bars) and model-predicted (blue lines) seroprevalence by age. The 95% credible intervals (95% CrIs) representing the 2.5th and 97.5th percentiles of the model’s posterior distributions are shown as the blue shading. Panel A shows the results from a model considering maternal antibody dynamics, and Panel B shows the model fit when neglecting maternal antibody dynamics.

year (95% CrI: 1.0–2.9), suggesting an average duration of maternal antibodies of 234 days or around 8 months, and (ii) an annual probability of infection of 0.18 (95% CrI: 0.15–0.20).

9.2 | Feedback Between Mothers and Their Babies

Whether a child has maternal immunity depends on the mother’s history of exposure to a pathogen and whether the mother retained transferable antibodies. This, in turn, depends on maternal age and the birth year of the mother. These can be embodied in the initial condition of the system:

$$M^b(0) = X^{b'}(b - b'), \quad S^b(0) = S^{b'}(b - b'), \quad X^b(0) = 0 \quad (59)$$

where $b' < b$ is the date/time when the mother was born, meaning $b - b' > 0$ is the age at which they give birth.

Accounting for maternal age at birth might be most relevant if there are substantial changes in seroprevalence around the child-bearing age, as well as considerable FOI in the early years of life for the mother.

10 | Discussion

This paper illustrates the uses of serocatalytic models to form epidemiological insights about pathogen transmission patterns. Whilst our paper surveyed a range of serocatalytic models, there are a host of other developments of these models that, for brevity, we omitted. This includes how to account for vaccination (e.g., [28, 29]), the impact of migration and immigration

(e.g., [30–32]), seasonality ([33]), cross-strain immunity ([34, 35]) and cross-reactivity of immunological assays (e.g., [36]). We also did not discuss how the basic reproduction value can be inferred by fitting serocatalytic models [37]. We aim to cover these topics in a follow-up article.

The interpretation and conclusions derived from serocatalytic models are generally contingent on the assumptions made about the underlying transmission dynamics and epidemiology. For example, using a model that neglects to account for pre-existing immunity due to vaccination in large swathes of the population would misrepresent serological dynamics [38, 39]. Nuisance factors, like the rate at which antibodies wane and become undetectable, must also be accounted for to correctly interpret serological data [40, 41].

The purpose of this manuscript is to focus on the mathematics of serocatalytic models rather than the steps required to fit these models to serological data. The types of serological data can, however, substantially affect what epidemiological knowledge can be recovered, and study designs for serosurveillance include (i) single cross-sectional sampling, (ii) repeated, periodic specimen collections of population-representative serum and (iii) longitudinal sampling of a study cohort (or longitudinal capture–recapture). Whilst cross-sectional serological data cannot typically differentiate between age-dependent versus time-dependent infection risk, longitudinal data may be able to separate these.

For epidemiological studies, serocatalytic models are typically fitted to serological data, and there are many nuances for doing so. We have developed an open-source R package, serofoi³, a fully tested software package for fitting serocatalytic models to data while incorporating best practices for model inference.

A recurring theme when considering inference for these models is that the conclusions drawn are highly contingent on the assumptions made. This is because often divergent theories of disease circulation can generate similar seroprevalence patterns. The epidemiological context is then crucial to ensure that serological data are correctly interpreted.

Acknowledgments

S.B. would also like to thank Merton College, University of Oxford, where she is the Peter J Braam Early Career Research Fellow in Global Wellbeing.

Ethics Statement

The authors have nothing to report.

Conflicts of Interest

The authors declare no conflicts of interest.

Data Availability Statement

Data sharing is not applicable to this article as no datasets were generated or analyzed during the current study.

Endnotes

¹ This neglects the effect of maternal antibodies, which we discuss in Section 6.

² In the stochastic individual-based version of this model.

³ <https://github.com/epiverse-trace/serofoi>.

References

1. J. Hay, I. Routledge, and S. Takahashi, “Serodynamics: A Review of Methods for Epidemiological Inference Using Serological Data,” *OSF Prepr* 2023.
2. H. Muench, “Derivation of Rates From Summation Data by the Catalytic Curve,” *Journal of the American Statistical Association* 29 (1934): 25–38.
3. Q. Clairon, M. Prague, D. Planas, et al., “Modeling the Kinetics of the Neutralizing Antibody Response Against SARS-CoV-2 Variants After Several Administrations of Bnt162b2,” *PLoS Computational Biology* 19 (2023): e1011282.
4. F. T. Cutts and M. Hanson, “Seroepidemiology: An Underused Tool for Designing and Monitoring Vaccination Programmes in Low-and Middle-Income Countries,” *Tropical Medicine & International Health* 21 (2016): 1086–1098.
5. C. J. E. Metcalf, J. Farrar, F. T. Cutts, et al., “Use of Serological Surveys to Generate Key Insights Into the Changing Global Landscape of Infectious Disease,” *Lancet* 388 (2016): 728–730.
6. K. Iversen, H. Bundgaard, R. B. Hasselbalch, et al., “Risk of COVID-19 in Health-Care Workers in Denmark: An Observational Cohort Study,” *Lancet Infectious Diseases* 20 (2020): 1401–1408.
7. R Core Team, *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, 2024).
8. W. M. Landau, “The Targets R Package: A Dynamic Make-Like Function-Oriented Pipeline Toolkit for Reproducibility and High-Performance Computing,” *Journal of Open Source Software* 6 (2021): 2959.
9. K. Ushey and H. Wickham, “renv: Project Environments,” 2024. R package version 1.0.11, <https://github.com/rstudio/renv>.

10. S. A. Vegt, L. Dai, I. Bouros, et al., “Learning Transmission Dynamics Modelling of COVID-19 Using Comomodels,” *Mathematical Biosciences* 349 (2022): 108824.
11. N. F. Brazeau, R. Verity, S. Jenks, et al., “Estimating the COVID-19 Infection Fatality Ratio Accounting for Seroreversion Using Statistical Modelling,” *Communications Medicine* 2 (2022): 54.
12. S. Kayoko, S. Y. Lau Max, N. M. Kraay Alicia, et al., “Estimating the Cumulative Incidence of SARS-CoV-2 Infection and the Infection Fatality Ratio in Light of Waning Antibodies,” *Epidemiology* 32, no. 4 (2021): 518–524.
13. K. Bartholomeeusen, D. Matthieu, D. A. LaBeaud, et al., “Chikungunya Fever,” *Nature Reviews Disease Primers* 9 (2023): 17.
14. J. K. Lim, V. Ridde, T. Agnandji Selidji, et al., “Seroepidemiological Reconstruction of Long-Term Chikungunya Virus Circulation in Burkina Faso and Gabon,” *Journal of Infectious Diseases* 227 (2023): 261–267.
15. A. Hviid, S. Rubin, and K. Mühlemann, “Mumps,” *Lancet* 371 (2008): 932–944.
16. C. P. Farrington, M. N. Kanaan, and N. J. Gay, “Estimation of the Basic Reproduction Number for Infectious Diseases From Age-Stratified Serological Survey Data,” *Journal of the Royal Statistical Society: Series C: Applied Statistics* 50 (2001): 251–292.
17. L. R. Baden, R. Kanapathipillai, E. W. Campion, S. Morrissey, E. J. Rubin, and J. M. Drazen, “Ebola—an ongoing crisis,” *New England Journal of Medicine* 371 (2014): 1458–1459.
18. D. Manno, P. Ayieko, D. Ishola, et al., “Ebola Virus Glycoprotein IgG Seroprevalence in Community Previously Affected by Ebola, Sierra Leone,” *Emerging Infectious Diseases* 28 (2022): 734–738.
19. A. Forna, P. Nouvellet, I. Dorigatti, and C. A. Donnelly, “Case Fatality Ratio Estimates for the 2013–2016 West African Ebola Epidemic: Application of Boosted Regression Trees for Imputation,” *Clinical Infectious Diseases* 70 (2020): 2476–2483.
20. M. Q. Ott, T. Bärnighausen, F. Tanser, M. N. Lurie, and M.-L. Newell, “Age-Gaps in Sexual Partnerships: Seeing Beyond ‘Sugar Daddies’,” *AIDS* 25 (2011): 861–863.
21. J. Todd, J. R. Glynn, M. Marston, et al., “Time From HIV Seroconversion to Death: A Collaborative Analysis of Eight Studies in Six Low and Middle-Income Countries Before Highly Active Antiretroviral Therapy,” *AIDS* 21 (2007): S55–S63.
22. M. Joël, G. Erofilii, T. Frank, B. Till, and N. Marie-Louise, “Modelling HIV Incidence and Survival From Age-Specific Seroprevalence After Antiretroviral Treatment Scale-Up in Rural South Africa,” *AIDS* 27 (2013): 2471–2479.
23. J. Todd, I. Cremin, N. McGrath, et al., “Reported Number of Sexual Partners: Comparison of Data From Four African Longitudinal Studies,” *Sexually Transmitted Infections* 85 (2009): i72–i80.
24. C. S. Grizer, K. Messacar, and J. J. Mattapallil, “Enterovirus-D68—a Reemerging Non-Polio Enterovirus That Causes Severe Respiratory and Neurological Disease in Children,” *Frontiers in Virology* 4 (2024): 1328457.
25. E. Kamau, H. Harvala, S. Blomqvist, et al., “Increase in Enterovirus D68 Infections in Young Children, United Kingdom, 2006–2016,” *Emerging Infectious Diseases* 25 (2019): 1200–1203.
26. M. Pons-Salort, B. Lambert, E. Kamau, et al., “Changes in Transmission of Enterovirus D68 (EV-D68) in England Inferred From Seroprevalence Data,” *eLife* 12 (2023): e76609.
27. D. Jorgensen, N. C. Grassly, and M. Pons-Salort, “Global Age-Stratified Seroprevalence of Enterovirus D68: A Systematic Literature Review,” *Lancet Microbe* 6, no. 1 (2024): 100938.

28. W. Wang, M. O’Driscoll, Q. Wang, S. Zhao, H. Salje, and H. Yu, “Dynamics of Measles Immunity From Birth and Following Vaccination,” *Nature Microbiology* 9 (2024): 1–1685.
29. F. Trentini, P. Poletti, S. Merler, and A. Melegaro, “Measles Immunity Gaps and the Progress Towards Elimination: A Multi-Country Modelling Analysis,” *Lancet Infectious Diseases* 17 (2017): 1089–1097.
30. H. E. Clapham, W. N. Chia, L. W. L. Tan, et al., “Contrasting SARS-CoV-2 Epidemics in Singapore: Cohort Studies in Migrant Workers and the General Population,” *International Journal of Infectious Diseases* 115 (2022): 72–78.
31. S. Simon, M. Amaku, and E. Massad, “Effects of Migration and Vaccination on the Spread and Control of Yellow Fever in Latin American Communities: A Mathematical Modelling Study,” *Lancet Planetary Health* 6 (2022): S7.
32. M. H. Bonds and P. Rohani, “Herd Immunity Acquired Indirectly From Interactions Between the Ecology of Infectious Diseases, Demography and Economics,” *Journal of the Royal Society Interface* 7 (2010): 541–547.
33. M. G. Johnsen, L. E. Christiansen, and K. Græsboell, “Seasonal Variation in the Transmission Rate of COVID-19 in a Temperate Climate Can Be Implemented in Epidemic Population Models by Using Daily Average Temperature as a Proxy for Seasonal Changes in Transmission Rate,” *Microbial Risk Analysis* 22 (2022): 100235.
34. R. C. Reiner, Jr., S. T. Stoddard, B. M. Forshey, et al., “Time-Varying, Serotype-Specific Force of Infection of Dengue Virus,” *Proceedings of the National Academy of Sciences of the United States of America* 111 (2014): E2694–E2702.
35. R. N. Thompson, E. Southall, Y. Daon, et al., “The Impact of Cross-Reactive Immunity on the Emergence of SARS-CoV-2 Variants,” *Frontiers in Immunology* 13 (2023): 1049458.
36. D. B. Larremore, B. K. Fosdick, K. M. Bubar, et al., “Estimating SARS-CoV-2 Seroprevalence and Epidemiological Parameters With Uncertainty From Serological Surveys,” *eLife* 10 (2021): e64206.
37. N. M. Ferguson, C. A. Donnelly, and R. M. Anderson, “Transmission Dynamics and Epidemiology of Dengue: Insights From Age-Stratified Sero-Prevalence Surveys,” *Philosophical Transactions of the Royal Society of London. Series B, Biological Sciences* 354 (1999): 757–768.
38. D. Nash, A. Srivastava, Y. Shen, et al., “Seroincidence of SARS-CoV-2 Infection Prior to and During the Rollout of Vaccines in a Community-Based Prospective Cohort of US Adults,” *Scientific Reports* 14 (2024): 644.
39. K. M. Wu and S. Riley, “Simulation-Guided Design of Serological Surveys of the Cumulative Incidence of Influenza Infection,” *BMC Infectious Diseases* 14 (2014): 1–9.
40. Vinh Dao Nguyen and F. Boni Maciej, “Statistical Identifiability and Sample Size Calculations for Serial Seroepidemiology,” *Epidemics* 12 (2015): 30–39.
41. N. Hens, M. Aerts, Z. Shkedy, H. Theeten, P. Van Damme, and P. Beutels, “Modelling Multisera Data: The Estimation of New Joint and Conditional Epidemiological Parameters,” *Statistics in Medicine* 27 (2008): 2651–2664.
42. A. Rohatgi, “WebPlotDigitizer: Version 4.5,” 2021.
43. Stan Development Team, “RStan: The R Interface to Stan,” 2024 R package version 2.32.6.

Appendix A

Muench’s Serosurvey Data for Yellow Fever and Methods for Their Analysis

Muench’s serological datasets for yellow fever are contained within [2]: For the two Brazil datasets, in their Chart 1; and for the Colombia dataset,

in their Table 1. The data from the charts were digitized using *webplotdigitizer* [42]. The chart-derived data had no information on the sample sizes, unlike the tabular data.

For all datasets, we used the method of maximum likelihood to estimate the FOIs. We used two different serocatalytic models to analyze the data and estimate the FOI.

Brazil Datasets

We assumed an age- and time-constant FOI, $\bar{\lambda}$, meaning the proportion of individuals seropositive by integer age A was given by:

$$X(A) = 1 - \exp(-\bar{\lambda}A) \quad (\text{A1})$$

Assuming counts of $x(A)$ seropositive individuals from a sample of $n(A)$ individuals of integer age A , this implies the probability of observing those data is given by:

$$\mathbb{P}(x(A)|n(A)) = \binom{n(A)}{x(A)} X(A)^{x(A)} (1 - X(A))^{n(A)-x(A)} \quad (\text{A2})$$

$$= \binom{n(A)}{x(A)} (1 - \exp(-\bar{\lambda}A))^{x(A)} \times \exp(-\bar{\lambda}a(n(A) - x(A))) \quad (\text{A3})$$

Considering a sample of n distinct age categories, where in age-category i (with individuals of age a_i) we observe x_i seropositive individuals out of n_i individuals tested, we can write down the likelihood as:

$$l = \prod_{i=1}^n \binom{n_i}{x_i} (1 - \exp(-\bar{\lambda}a_i))^{x_i} \exp(-\bar{\lambda}a_i(n_i - x_i)) \quad (\text{A4})$$

From Equation (A4), we can write down the log-likelihood up to a constant:

$$\mathcal{L} \sim \sum_{i=1}^n x_i \log(1 - \exp(-\bar{\lambda}a_i)) - (n_i - x_i)\bar{\lambda}a_i \quad (\text{A5})$$

Taking the derivative of Equation (A5) and setting this to zero, we have the condition satisfied by the maximum likelihood estimator:

$$\frac{\partial \mathcal{L}}{\partial \bar{\lambda}} = \sum_{i=1}^n \left(\frac{x_i a_i \exp(-\bar{\lambda}a_i)}{1 - \exp(-\bar{\lambda}a_i)} - (n_i - x_i)a_i \right) = 0 \quad (\text{A6})$$

Equation (A6) contains n_i , which is unknown for the Brazil datasets. But the Brazil data contained only the proportions seropositive, p_i , and we assume that $x_i = p_i n_i$; substituting this into Equation (A6) leaves a factor n_i in all terms which can be divided through by.

The solution to Equation (A6) can only be determined numerically, and we used R’s *optim* function to locate the maximum likelihood estimates for each of the two datasets.

Colombia Dataset

We assumed that anyone alive at the time of the 1929 yellow fever epidemic was equally likely to be exposed to the virus and experienced an FOI, λ . We assume the serosurvey was conducted τ years after the epidemic. This results in a proportion seropositive given by:

$$X(a) = \begin{cases} 1 - \exp(-\lambda), & \text{if } a > \tau, \\ 0, & \text{otherwise} \end{cases} \quad (\text{A7})$$

The youngest age group in the dataset is 5–9, and since the study was published in 1934, we assume that everyone could have been exposed to the virus (i.e., $a > \tau$ for everyone in the sample).

Through similar logic to before, we can determine a condition for the maximum likelihood estimator, $\hat{\lambda}$, although this time in closed form:

$$\hat{\lambda} = -\log\left(1 - \frac{\sum_{i=1}^n x_i}{\sum_{i=1}^n n_i}\right) \quad (\text{A8})$$

Glossary S1: **Mathematical notations used in this paper.**

- t the continuous calendar date/time.
- T the integer calendar year (e.g., 2024).
- b individuals' continuous calendar date/time of birth.
- B individuals' calendar year of birth (e.g., 1985).
- S^b the proportion of individuals born at calendar time b who are **seronegative** at time t .
- X^b the proportion of individuals born at calendar time b who are **seropositive** at time t .
- $\bar{\lambda}$ a constant FOI, that is, independent of age and time, where $\bar{\lambda} > 0$.
- $\bar{\lambda}_T$ an annual FOI, assumed to potentially vary across years but be constant within calendar year T .
- $\bar{\lambda}_A$ an age-varying FOI, assumed to be constant within a calendar year of age, A , but varying across ages.
- $a^b = t - b$, for any $t \geq \tau$ the **age in continuous time units** of individuals who are born at calendar time b , evaluated at time t . Here, b is fixed since this age corresponds to a single birth cohort.
- $a_t = t - b$, for any $t \geq \tau$ the **age in continuous time units** of individuals who are born at calendar time b , evaluated at time t . Here, t is fixed and b varies since we consider a cross-section of the population at time t .
- $A^B = T - B$ the **age in integer year units** of individuals who are born in year B , evaluated in year T . Here, B is fixed since this age corresponds to a single birth cohort.
- $A_T = T - B$ the **age in integer year units** of individuals who are born in year B , evaluated in year T . Here, T is fixed and B varies since we consider a cross-section of the population in year T .
- μ a constant rate of **seroreversion**, where $\mu > 0$.
- γ a constant rate of waning of **maternal antibodies**, where $\gamma > 0$.
- ϵ a constant **death** rate due to current or past infection, where $\epsilon > 0$.
- ρ represents the proportion of infections which are severe (resulting in death); alternatively, known as the **infection fatality ratio**, where $0 \leq \rho \leq 1$.

Transmission Dynamics Model Including Waning Immunity

To generate Figure 2, we simulated from a transmission dynamics model of the following form:

$$\begin{aligned} \frac{dS(t)}{dt} &= -\beta(t)S(t)I(t) + \delta R(t), \\ \frac{dI(t)}{dt} &= \beta(t)S(t)I(t) - \gamma I(t), \\ \frac{dR(t)}{dt} &= \gamma I(t) - \delta R(t) \end{aligned} \quad (\text{A9})$$

where $\delta = 0.002$ per year is the rate at which recovered individuals become susceptible again (i.e., their immunity wanes), and $\gamma = 0.05$ per

year is the rate at which infected individuals recover. We assume that the rate of infection, $\beta(t)$, is piecewise-constant and given by:

$$\beta(t) = \begin{cases} 0.06, & \text{if } t \leq 365, \\ 0.1, & \text{if } 365 < t \leq 2 \times 365, \\ 0.1, & \text{if } 2 \times 365 < t \leq 3 \times 365, \\ 0.05, & \text{if } 3 \times 365 < t \end{cases} \quad (\text{A10})$$

The initial conditions for the system given by Equation (A9) were: $S(0) = 0.99$, $I(0) = 0.01$, $R(0) = 0$, and the system was solved numerically using Mathematica's inbuilt NDSolve method.

Bayesian Inference for Serocatalytic Models

Throughout this paper, we show that serocatalytic models can be fitted to age-structured seroprevalence data to yield estimates of FOIs and other quantities.

Here, we provide an overview of the approach used for inference in this paper. In each of the real data examples, our dataset consisted of serosurveys conducted at a specific time (which we assume is at the start of the integer year T). The serosurvey results can be summarized by a series of pairs (x_A, n_A) , where x_A denotes the count of seropositive individuals of integer age A within those sampled, n_A . Throughout, we assume that the proportion seropositive follows:

$$x_A \stackrel{\text{iid}}{\sim} \text{binomial}(n_A, X_A(T)) \quad (\text{A11})$$

where $X_A(T)$ is the solution to a serocatalytic model and represents the seropositivity for individuals of age A at the start of calendar year T in the population where the survey was conducted.

In some of the analyses, we assumed that the FOIs are piecewise-constant; in others, we assumed that the FOI followed a parametric curve and inferred the parameters that specified that curve. In Table A1, we describe the assumptions made for each of the model fits.

We used a Bayesian framework for inference, which requires prior distributions to be specified for all unknown model parameters. In Table A1, we provide the prior distributions used in each of the analyses. These were chosen to allow a wide range of model solutions.

To fit the models, we used Markov chain Monte Carlo (MCMC) sampling through Stan's default NUTS algorithm [43]. For each model fit, we used ≥ 3000 MCMC iterations per chain across four Markov chains. Model convergence was diagnosed by $\hat{R} < 1.1$ and effective sample sizes (bulk and tail) above 200 for all parameters.

TABLE A1 | Priors used in model fitting.

Model	Priors	Corresponding section
Time-varying FOI, λ_i , where i is the index of an FOI piece	$\log(\lambda_{i>1}) \sim \text{Student-t}(v, \log(\lambda_{i-1}), \sigma)$, $\log(\lambda_{i=1}) \sim \text{normal}(-3, 1)$, $\sigma \sim \text{Cauchy}^+(0, 1)$, $v \sim \text{Cauchy}^+(0, 1)$.	Section 5.2
Age-varying FOI, λ_i , where i is the index of an FOI piece	$\log(\lambda_{i>1}) \sim \text{normal}(\log(\lambda_{i-1}), \sigma)$, $\log(\lambda_{i=1}) \sim \text{normal}(-3, 1)$, $\sigma \sim \text{Cauchy}^+(0, 1)$.	Section 6.3
Elevated death rate due to infection	$\lambda \sim \text{exponential}(1)$, $\kappa \sim \text{normal}^+(0.04, 0.02)$.	Section 7.2
Time- and age-dependent FOI	$a \sim \text{Cauchy}^+(0, 1)$, $b \sim \text{Cauchy}^+(0, 1)$, $c \sim \text{Cauchy}^+(0, 1)$	Section 8.2
Maternal antibodies	$\lambda \sim \text{exponential}(1)$, $\gamma \sim \text{Cauchy}^+(0, 1)$.	Section 9.1

Note: The right column lists the respective section in the paper where the model is fitted to a real serosurvey dataset. The parameter ϵ in the *elevated death rate due to infection* model is given by $\epsilon = 1/\kappa$, where $\kappa > 0$ represents the typical time until death. The age-component of the FOI in the *time- and age-dependent* model is parametrically modeled as a function of a , b , and c parameters shown in this table: $\lambda = c \times \text{gamma}(\text{age}|a, b)$, where $\text{gamma}(\text{age}|a, b)$ is a gamma distribution probability density function evaluated at age, which has mean a/b ; the time component in the FOI is given a uniform prior between 0 and 1. The λ_i in Section 5.2 and λ_a in Section 6.3 (Panels B and D in Figure 7) were divided into smaller groups, or “chunks”, and the FOI for each chunk was estimated. See https://github.com/ekamau/serocatalytic_models for further details.